

Oliver Freeman

**A Comparison of Any-Angle Pathfinding
Algorithms for Virtual Agents**

Computer Science: Part II

Clare College

January 24, 2014

Proforma

Name: **Oliver Freeman**
College: **Clare College**
Project Title: **A Comparison of Any-Angle Pathfinding Algorithms for Virtual Agents**
Examination: **Computer Science Tripos: Part II, May 2014**
Word Count: **????¹**
Project Originator: **Oliver Freeman**
Supervisor: **Dr R. J. Gibbens**

Original Aims of the Project

Work Completed

Special Difficulties

None

Original
Aims
of the
Project
...

Work
com-
pleted
...

Declaration

I, Oliver Freeman of Clare College, being a candidate for Part II of the Computer Science Tripos, hereby declare that this dissertation and the work described in it are my own work, unaided except as may be specified below, and that the dissertation does not contain material that has already been used to any substantial extent for a comparable purpose.

Signed [signature]

Date [date]

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Related work	2
1.3	Project goals	2
2	Preparation	3
2.1	Introduction to any-angle pathfinding	3
2.1.1	Definitions	3
2.1.2	The any-angle pathfinding problem	4
2.1.3	Discretization	4
2.2	Any-angle pathfinding algorithms	5
2.2.1	Dijkstra’s shortest paths	6
2.2.2	A*	6
2.2.3	A* with post-smoothing	7
2.2.4	Basic θ^*	8
2.2.5	Lazy θ^*	8
2.2.6	Block A*	10
2.3	Requirements analysis	12
2.3.1	Testing simulator	12
2.3.2	Map generation	12
2.3.3	Algorithm implementation	13
2.3.4	Data gathering	13
2.4	Testing	13
2.5	Design model	13
2.6	Languages and tools	14
3	Implementation	15
3.1	Map generation	16
3.2	Simulation	17
3.2.1	Graph Generation	17
3.2.2	Algorithm Data	18
3.3	Algorithms	19
3.3.1	Dijkstra’s shortest paths	22
3.3.2	A*	22
3.3.3	A* with post-smoothing	22
3.3.4	Basic θ^* and Lazy θ^*	22
3.3.5	Block A*	22
3.4	User Interface	26
3.5	Testing	26
3.6	Data extraction	27

4	Evaluation	29
5	Conclusion	30

List of Figures

1.1	Shortest paths through maps and graphs	1
2.1	A node and its eight neighbours	4
2.2	Different graph representations for a given map	5
2.3	A valid path for an agent modelled as a point	5
2.4	A* with post-smoothing	7
2.5	Block A*: map of 9×9 cells is split into 9 blocks of 3×3 cells	11
3.1	Flow of the user interface and simulator	16
3.2	Two iterations of GenerateMap with $R=4$ and $D=2$	16
3.3	Code snippet showing Node creation for grid-based graphs	18
3.4	Line of sight algorithm	19
3.5	Composition of MapInstance	20
3.6	Inheritance structure of algorithms	21
3.7	Extract of the LDDb for block size of 3×3 - array of <code>HashTable<PairOfCoords,Pair<double,List<</code>	
3.8	Encoding scheme for <code>PairOfCoords</code> in the compressed LDDb	24
3.9	Encoding scheme for <code>List<Coordinate></code> in the compressed LDDb	24
3.10	User interface in map creation mode	26
3.11	User interface displaying paths for a generated map	27
3.12	Extract from a CSV file exported by DataExtract	27
3.13	Code snippet from DataExtraction showing automation of Graph creation	28

Acknowledgements

Acknowledgements
including
citations?
...

Chapter 1

Introduction

1.1 Motivation

Finding short and realistic-looking paths through maps with arbitrarily placed obstacles is one of the central problems in artificial intelligence for games and robotics. Figure 1.1 (a) shows a grid-based map, consisting of free cells and blocked cells. Figure 1.1 (c) shows the shortest path through the map between the *start* and the *goal*.

However, path-finding algorithms operate on graphs. A graph representation of the map is shown in figure 1.1 (b) — where each node in the graph represents a coordinate on the map and a character can travel directly between any two nodes connected by an edge. Figure 1.1 (d) shows the optimum path through the graph — however, this is clearly longer than the optimum path through the map.

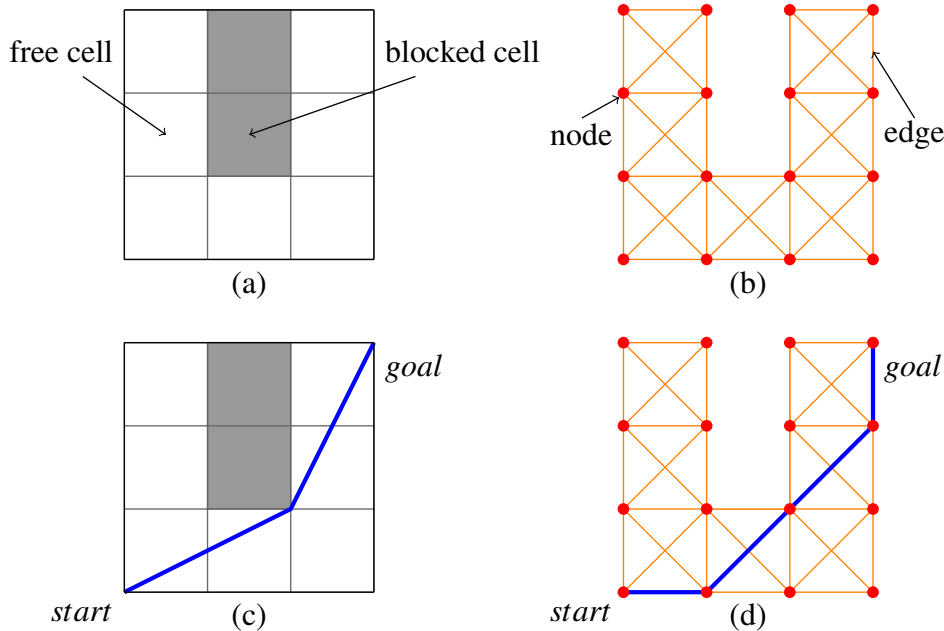


Figure 1.1: Shortest path through a map vs. shortest path through the graph representing the map

If the graph representation were changed then a more realistic path could be achieved, but in general such a graph would require many more nodes and edges which would cause exponentially increasing memory requirements. This dissertation will investigate various algorithms that aim to find optimal paths through maps while using the space-efficient grid-based graphs show in figure 1.1(b).

1.2 Related work

Related
work
...

1.3 Project goals

Project
goals
...

Chapter 2

Preparation

2.1 Introduction to any-angle pathfinding

In this section I will introduce some key terms, define the any-angle pathfinding problem, and describe some decisions and methods to solve the problem.

2.1.1 Definitions

Agent

An entity that can move from one location in the environment to another. The agent cannot move to or through a location that is blocked by an obstacle.

Map

A representation of an environment as a two dimensional plane \mathbb{R}^2 . For simplicity, this dissertation will focus on grid-based maps, where a map is divided into N^2 cells of equal size.

Cell

The logical unit of a map. Each cell may either be free or blocked. A blocked cell represents an obstacle that completely fills the cell.

Line of Sight

Exists between two coordinates on a map if a straight line between those two coordinates does not intersect any blocked cells.

Path

An ordered list of coordinates $(c_{start}, c_1, \dots, c_{goal})$. For the path to be valid, a line of sight must exist between c_{start} and c_1 , c_1 and c_2 , ... c_{n-1} and c_{goal} .

Graph

A finite mathematical construct, consisting of nodes and undirected edges. A graph is a discrete representation of the \mathbb{R}^2 space of the map.

Node

Represents a single *coordinate* on the map.

Edge

An unordered pair of nodes. The existence of edge (n_1, n_2) denotes that an agent can travel in a straight line between the coordinates represented by n_1 and n_2 . Each edge has an associated *weight*, which represents the Euclidean distance between the coordinates represented by the two nodes.

Neighbour

A node n' is the neighbour of a node n if there exists an edge (n, n') .

2.1.2 The any-angle pathfinding problem

The problem is to compute optimal or near-optimal paths, if they exist, through maps. The optimality is based on the total Euclidean distance and cumulative angle turned by an agent following the path.

The problem is split into two sub-problems:

- *Discretization* - a graph is created that represents that map in some way
- *Pathfinding* - an algorithm is run on the graph to produce a path

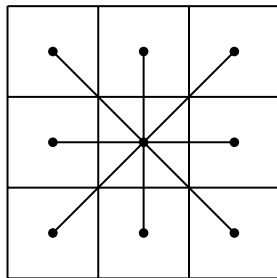
These two sub-problems will now be introduced:

2.1.3 Discretization

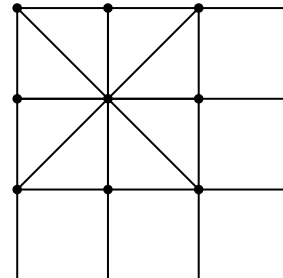
We must choose how to represent a continuous space, grid-based map with a graph consisting of a finite number of nodes and edges.

Grid-based graph

A common representation is to use a grid-based graph: there is one node per cell, and there are edges between nodes that represent locations that can be reached in a 1 unit step at 0° , 90° , 180° or 270° bearing, or a $\sqrt{2}$ step at 45° , 135° , 225° or 315° bearing. If the location represented by one node cannot be reached from the location represented by another node (because there is a blocked cell in between), then there will be no edge between the two nodes. We are now presented with a choice as to where to place the nodes:



(a) Nodes in the centres of cells



(b) Nodes on the corners of cells

Figure 2.1: A node and its eight neighbours

The choice is essentially a stylistic one for most of the algorithms, but placing nodes in the middle of cells places stringent restrictions on the paths produced by Block A*, so I will place my nodes on the corners of cells.

Visibility graph

This is an alternative graph representation. Nodes exist at the start and end coordinates, and otherwise only on the corner of cells for which exactly three of the four surrounding cells are free. An edge exists between any two nodes that have a line of sight.

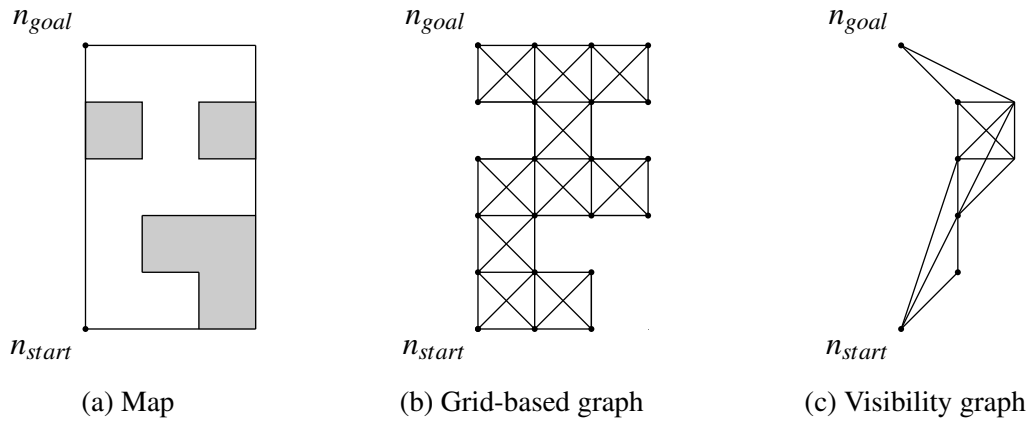


Figure 2.2: Different graph representations for a given map

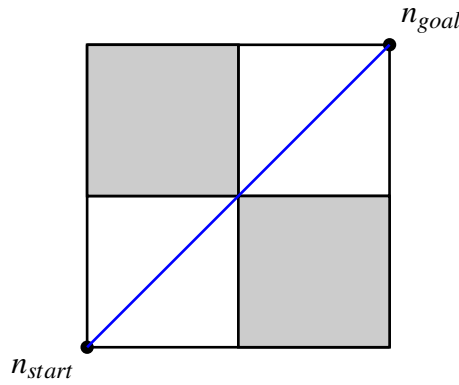


Figure 2.3: A valid path for an agent modelled as a point

The optimal path through a visibility graph will correspond to the optimal path through a map, whereas the optimal path through a grid-based graph may not. However, grid-based graphs are generally accepted as preferable since, for a map of N^2 cells, a grid-based graph will have $O(N^2)$ edges, whereas a visibility graph has $O(N^4)$ edges, which can be unfeasible for large or high resolution maps. Therefore, this investigation will focus predominantly on grid-based graphs.

Agent size

I will adopt the standard practice in any-angle pathfinding research by modelling the agent as a dimensionless point. Therefore, the path shown in Figure 2.3, which features a 'diagonal blockage', is a valid path. It is possible to ensure that such paths are never taken by adding extra checks in the graph creation and line of sight algorithms, but this would distract from the core investigation and be an unnecessary deviation from the accepted standard.

2.2 Any-angle pathfinding algorithms

In this section, I will introduce the four any-angle pathfinding algorithms that I will be investigating: A* with post smoothing, Theta*, Lazy Theta* and Block A*. Firstly, I will introduce an additional two comparatively basic algorithms to aid in the explanation of the more complex algorithms.

2.2.1 Dijkstra's shortest paths

Most of the algorithms that I will study will be derivatives of the well known A* graph traversal algorithm, which itself is a derivative of Dijkstra's famous shortest-path algorithm.

Dijkstra finds the shortest path from a given node n_{start} to all other nodes in a graph - though it can be prematurely terminated when the goal is reached. Each node n has a:

- *g-value* - represents the length of the shortest path found so far from n_{start} to n . Initialised to ∞ .
- *parent* - the previous node in the shortest path found so far from n_{start} to n

The algorithm processes each node once only, selecting it from *openSet*, a priority queue that prioritises low *g*-values. When a node is selected, its *g*-value is the length of the shortest path to it from n_{start} . A closed set is used to ensure that no node is expanded twice, as this would incur unnecessary work. When n_{goal} is processed, the algorithm terminates. Dijkstra's shortest-path algorithm is optimal and complete.

Algorithm 1: DIJKSTRA

```

def Dijkstra( $G, n_{start}, n_{goal}$ )
1   $openSet \leftarrow \emptyset$ 
2   $closedSet \leftarrow \emptyset$ 
3   $n_{start}.g \leftarrow 0$ 
4   $openSet.add(n_{start})$ 
5  while  $openSet \neq \emptyset$  do
6       $n_{curr} \leftarrow openSet.pop()$ 
7       $closedSet.add(n_{curr})$ 
8      if  $n_{curr} = n_{goal}$  then
9          return  $n_{goal}$ 
10     foreach  $n_{neigh}$  of  $n_{curr}$  do
11         if  $closedSet.contains(n_{neigh}) = false$  then
12             if  $Update(n_{neigh}) = true$  then
13                 if  $openSet.contains(n_{neigh}) = false$  then
14                      $openSet.add(n_{neigh})$ 
15 return  $\emptyset$ 

def Update( $n_{neigh}$ )
1  if  $n_{curr}.g + euclidean(n_{curr}, n_{neigh}) < n_{neigh}.g$  then
2       $n_{neigh}.g = n_{curr}.g + euclidean(n_{curr}, n_{neigh})$ 
3       $n_{neigh}.parent = n_{curr}$ 
4      return true
5  else
6      return false

```

If a path exists, the algorithm returns n_{goal} . The shortest path from n_{start} to n_{goal} can be then be traced: ($n_{goal}, n_{goal}.parent, n_{goal}.parent.parent, \dots, n_{start}$)

2.2.2 A*

A* is based on Dijkstra's shortest-path algorithm, but uses a heuristic h to reduce the number of nodes expanded, and hence often finds the optimal route to n_{goal} in less time. Where Dijkstra preferentially

expands nodes with low g -values, A* preferentially expands nodes with low f -values, where $f(n) = g(n) + h(n)$:

- $g(n)$ - calculated in the same way as in Dijkstra's algorithm
- $h(n)$ - Euclidean distance between n and n_{goal} - a cheaply computable monotonic estimate of the actual distance that will need to be traversed.

The monotonicity of Euclidean distance as a heuristic ensures that A* is optimal and complete.

The pseudocode for A* differs only from Dijkstra in the update subroutine, where the h -score must also be initialised.

Algorithm 2: Update from A*

```

def Update( $n_{neigh}$ )
1  if  $n_{curr}.g + euclidean(n_{curr}, n_{neigh}) < n_{neigh}.g$  then
2       $n_{neigh}.g \leftarrow n_{curr}.g + euclidean(n_{curr}, n_{neigh})$ 
3       $n_{neigh}.f \leftarrow euclidean(n_{neigh}, n_{goal})$ 
4       $n_{neigh}.parent = n_{curr}$ 
5      return true
6  else
7      return false

```

2.2.3 A* with post-smoothing

A post processing step is run on the path returned by A*, where any unnecessary deviations are cut out by changing the parents of nodes of the path:

The algorithm starts with the node returned by A*: $n = n_{goal}$. The algorithm performs a line of sight test from n to n 's parent's parent. If the test passes, n will have its parent set to be its parent's parent. This test and parent resetting is repeated until the line of sight test fails. Once the test has failed, we choose a new n , which will be the old n 's parent's parent - i.e. the node on which the line of sight test failed. We repeat this process until we reach the start node, at which point the algorithm terminates.

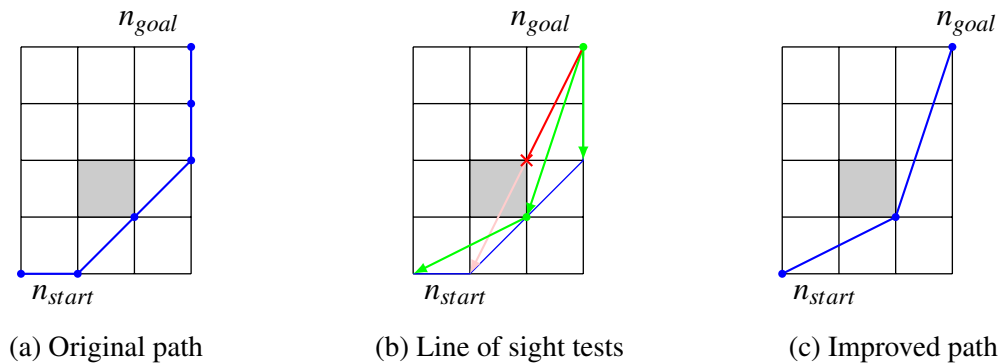


Figure 2.4: A* with post-smoothing

Algorithm 3: PostSmoothing from A* WITH POST-SMOOTHING

```

def PostSmoothing( $n_{start}, n_{goal}$ )
1   $n_{curr} \leftarrow n_{goal}$ 
2   $n_{next} \leftarrow n_{goal}.parent.parent$ 
3  if  $n_{next} = \emptyset$  then
4    return
5  while true do
6    while LineOfSight( $n_{curr}, n_{next}$ ) do
7       $n_{curr}.parent \leftarrow n_{next}$ 
8       $n_{next} \leftarrow n_{next}.parent$ 
9      if  $n_{next} = n_{start}$  then
10       return
11      $n_{curr} \leftarrow n_{next}$ 
12     if  $n_{curr}.parent = n_{start}$  then
13       return
14      $n_{next} \leftarrow n_{next}.parent.parent$ 

```

2.2.4 Basic θ^*

Where A* with post-smoothing does an explicit smoothing step after finding the basic path, Basic θ^* interleaves smoothing with exploration by attempting to re-parent each node that is expanded with the parent of its parent. This reparenting occurs if there exists a line of sight between the node being expanded and the parent of its parent.

Algorithm 4: Update from θ^*

```

def Update( $n_{neigh}$ )
1  if LineOfSight( $n_{neigh}, n_{curr}.parent$ ) = true then
2    if  $n_{curr}.parent.g + euclidean(n_{curr}.parent, n_{neigh}) < n_{neigh}.g$  then
3       $n_{neigh}.g \leftarrow n_{neigh}.parent.g + euclidean(n_{curr}.parent, n_{neigh})$ 
4       $n_{neigh}.f \leftarrow euclidean(n_{neigh}, n_{goal})$ 
5       $n_{neigh}.parent \leftarrow n_{curr}.parent$ 
6      return true
7    else
8      return false
9  else
10   if  $n_{curr}.g + euclidean(n_{curr}, n_{neigh}) < n_{neigh}.g$  then
11      $n_{neigh}.g \leftarrow n_{curr}.g + euclidean(n_{curr}, n_{neigh})$ 
12      $n_{neigh}.f \leftarrow euclidean(n_{neigh}, n_{goal})$ 
13      $n_{neigh}.parent \leftarrow n_{curr}$ 
14     return true
15   else
16     return false

```

2.2.5 Lazy θ^*

Lazy θ^* reduces the number of line of sight tests performed by Basic θ^* . Basic θ^* will do a line of sight on every node that is visited (i.e. every neighbour of every expanded node) - however, if the

node is visited but not expanded then this line of sight test will have been unnecessary.

Instead, when visiting a node n , Lazy θ^* assumes that there is a line of sight between each neighbour n' and $n.parent$, and updates the g-value and parent of each n' accordingly. The algorithm only checks to see if that line of sight actually exists if the n' is ever expanded, by calling `Initialise` on n' when it is popped off `openSet`. If the line of sight doesn't pass, `Initialise` alters the g-value of n' to reflect this by taking as its new parent the expanded neighbour n'' or n' which minimises $g(n') + distance(n', n'')$.

Algorithm 5: Initialise and Update from LAZY θ^*

```

def Initialise( $n_{curr}$ )
1  if LineOfSight( $n_{curr}, n_{curr}.parent$ ) = false then
2       $newParent \leftarrow \underset{n' \in expandedNeigh(n_{curr})}{\operatorname{argmin}} (n'.g + distance(n', n_{curr}))$ 
3       $n_{curr}.parent \leftarrow n'$ 
4       $n_{curr}.g \leftarrow n'.g + distance(n', n_{curr})$ 
def Update( $n_{neigh}$ )
    // assume line of sight test passes
5  if  $n_{curr}.parent.g + euclidean(n_{curr}.parent, n_{neigh}) < n_{neigh}.g$  then
6       $n_{neigh}.g \leftarrow n_{neigh}.parent.g + euclidean(n_{curr}.parent, n_{neigh})$ 
7       $n_{neigh}.f \leftarrow euclidean(n_{neigh}, n_{goal})$ 
8       $n_{neigh}.parent \leftarrow n_{curr}.parent$ 
9      return true
10 else
11     return false

```

2.2.6 Block A*

Block A* is by far the most complicated algorithm that I will investigate. It was published in 2011, and is at the very cutting edge of any-angle path-finding algorithmic research.

Block A* is a variant of A* that uses the concept of an LDDDB ('Local Distance Database') - a database that holds the distances and inflection points of the optimum routes through small n by n submaps, known as blocks. The size of the blocks in the LDDDB is optional - the size chosen will affect the size of the database, the lookup time of the database, the speed of the algorithm and the optimality of the algorithm.

Block A* deals with blocks in a similar way to that in which A* deals with nodes: blocks have neighbours and priority-values (called 'heap-values', like A*'s f-values), and can be expanded, put in a priority queue and popped off the priority queue. However, note that blocks can be re-expanded if a better route is found - this does not occur in A*.¹ The majority of the work of Block A* is performed by only considering nodes on the boundaries of blocks - since the LDDDB can be used to quickly give costs for traversing the inside of the nodes. Therefore, Block A* can quickly skip across the map, block by block. The blocks that contain n_{start} and n_{goal} are special cases, as there is no guarantee that these nodes lie on the boundary of their respective blocks - and they may even lie in the same block. Therefore, the LDDDB cannot be used for these blocks.

A block's heap-value is the lowest f-value of a boundary node that has been updated since the block was last expanded, and is reset to ∞ when the block is removed from the *openSet*. In a similar way to A*, the basic cycle to Block A* is that the block with the lowest heap-value will be removed from the *openSet* and expanded. To expand the block:

- a list Y is made of all boundary nodes of the block that have been updated since the block was last expanded.
- For each side of the block that isn't on the edge of the map (see Figure 2.4):
 - a list $listX$ is made of each node on that side
 - a corresponding list $listX'$ is made which contains, for each x_i in $listX$, an x'_i which is the corresponding node to x_i in the neighbour block to that side.
 - the g-value of x_i is updated if there is a $g(y) + LDDDB(y, x_i)$ that is smaller than $g(x_i)$. We call y the *ingress* coordinate and x_i the *egress* coordinate, as y is the point of entry to the block and x_i is the proposed point of exit from the block.
 - if x_i is updated, so is x'_i - to $g(x'_i) + distance(x_i, x'_i)$ - which is 0 if we use nodes on corners or 1 if we use nodes in centre of cells.
 - if the smallest f-value of an updated x' is smaller than the neighbour block's current heap-value, we update it and insert the neighbour block into the *openSet* if it is not already there.

Special cases - where, in the general case, the LDDDB is not sufficient:

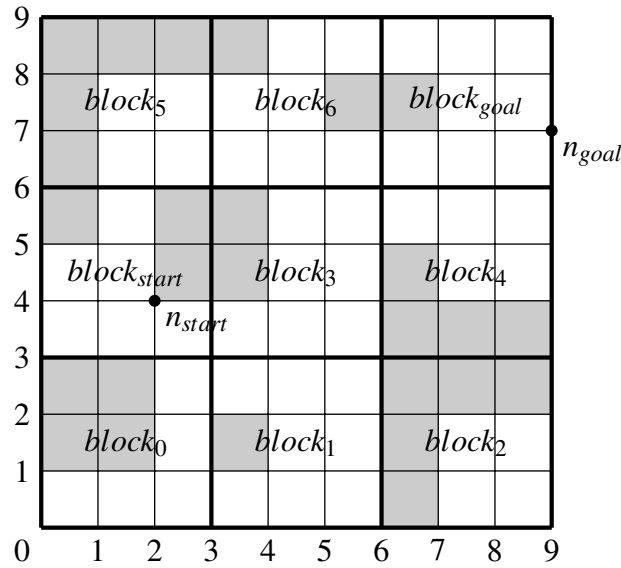
- $block_{start}$: the call to `init(start)` creates and initialises $block_{start}$. The optimum path from n_{start} to each boundary node x is computed, and the length of this path is used to set x 's g-value.

¹This subtlety caused me such confusion that I eventually emailed Peter Yapp, the author of the Block A* paper, for clarification. He explained the difference, and told me that he had included an explanation of this specific point when he gave a presentation [citation] on Block A* as the confusion is not uncommon.

- $block_{goal}$: the call to $init(goal)$ creates and initialises $block_{goal}$. The optimum path from n_{goal} to each boundary node x is computed, and the length of this path is used to set x 's h-value.²
- $block_{start} = block_{goal}$: the optimum path from n_{start} to n_{goal} is calculated directly, and returned.

When $block_{goal}$ is evaluated, its h-value has already been calculated by finding the actual shortest path from each boundary node to n_{goal} (see Special cases above), as opposed to using an estimate. Therefore, if the heap-value of $block_{goal}$ is less than any other blocks in $openSet$ then there are no possible shorter paths to $goal_{block}$, so the algorithm terminates.

The pseudocode for Block A* can be found in Appendix ?



When expanding $block_4$:

$listX =$
 $(6,3),(7,3),(8,3),(9,3),(6,4),(6,5),(6,6),(7,6),(8,6),(9,6)$ from $block_4$

$listX' =$
 $(6,3),(7,3),(8,3),(9,3)$ from $block_2$,
 $(6,3)$ from $block_1$,
 $(6,3),(6,4),(6,5),(6,6)$ from $block_3$,
 $(6,6)$ from $block_6$,
 $(6,6),(7,6),(8,6),(9,6)$ from $block_{goal}$

Figure 2.5: Block A*: map of 9×9 cells is split into 9 blocks of 3×3 cells

²h-value, not g-value. The difference between the treatment of $block_{start}$ and $block_{goal}$ is important

2.3 Requirements analysis

As specified in the ‘Work to be done’ section of my Proposal (see Appendix E), my project was divided into four sections. This section outlines the functional and non-functional requirements for the system, and their relative priorities using the MoSCoW system.

M - Must; **S** - Should; **C** - Could; **W** - Won’t

2.3.1 Testing simulator

ID	Functional requirement	Priority
1	The system shall load one of a collection of maps from the generator	M
2	The system shall load one of a collection of maps from a saved file	S
3	The system shall create a grid-graph from a given map	M
4	The system shall create a visibility graph from a given map	C
5	The system shall run one of a collection of any-angle path-finding algorithms on a graph and collect data such as the path-length and the length of computation	M
6	The system shall display a visual representation of the current map and the paths found by any algorithms that have been run on it	M
7	The system shall display the numeric statistics for each path for the current map	M
ID	Non-functional requirement	Priority
1	The system shall be designed in a modular way to allow easy extension for new algorithms	S

2.3.2 Map generation

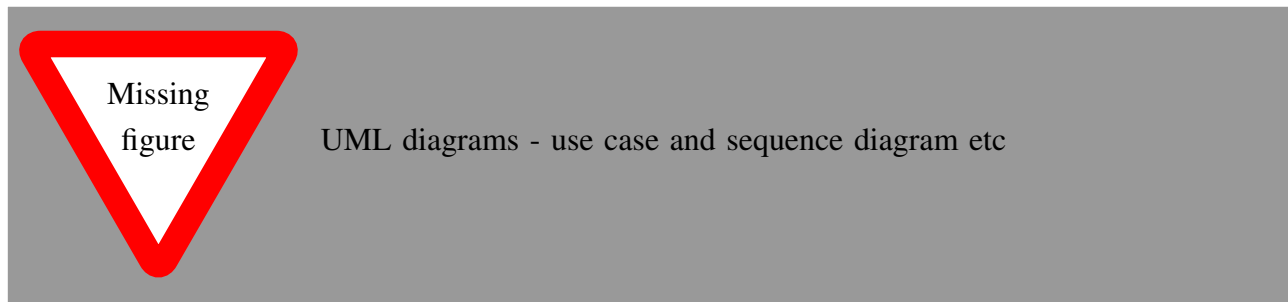
ID	Functional requirement	Priority
1	The system shall generate pseudo-random maps of a given resolution, coverage percentage and clustering	M
2	The system shall allow maps to be saved so that multiple tests can be run on the same map suite	M
3	The system shall allow maps to be created with an interactive map editor	C
ID	Non-functional requirement	Priority
1	The system shall generate maps of the highest resolution in under 2 seconds	S

2.3.3 Algorithm implementation

ID	Functional requirement	Priority
1	The system shall correctly implement each of the chosen algorithms. If a path exists, the path and numerical statistics will be returned. If no path exists, this will be returned	M
2	The system shall allow arbitrary start and end coordinates for any map	S
ID	Non-functional requirement	Priority
1	The system shall be designed in a modular way to allow easy extension for new algorithms	S

2.3.4 Data gathering

ID	Functional requirement	Priority
1	The system shall write statistics for an arbitrary set of specified algorithms on an arbitrary set of specified maps and write the results to a CSV file	M
ID	Non-functional requirement	Priority
1	The system shall be designed with a clear API that enables quick and easy data gathering.	M



2.4 Testing

2.5 Design model

Having completed the preparation phase of the project, I refined the plan for the implementation phase from that presented in the Project Proposal (see Appendix ?). An incremental model of implementation was adopted, with new modules being developed and tested separately before being integrated into the work program. The milestones of the project were:

Milestone 1 - Maps of arbitrary size, coverage and clustering can be created and printed to system output.

Milestone 2 - Arbitrary maps can be converted to graph, and A* can be run on these maps. A visual representation of the path can be printed to system output.

Milestone 3 - Basic UI is built, including all functionality from *Milestone 2*. Basic path statistics are displayed.

Milestone 4 - Map saving, map loading and map creation functionality are present. This will facilitate debugging edge cases for more complex algorithms.

Milestone 5 - Line of sight, A* with post-smoothing, θ^* and Lazy θ^* are implemented.

Milestone 6 - Block A* is implemented.

Milestone 7 - Data extraction scripts are implemented.

2.6 Languages and tools

Programming language

Java - provides abstraction and class hierarchy to enable development of modular, extensible code.

Libraries

Swing - for graphical user interface design and *CSVWriter* for data export.

Integrated development environment

Eclipse - allows rapid development through integrated testing, refactoring and version control tools.

Statistical analysis and visualisation

R - an open-source statistical package, due to its flexibility and extensibility.

Backup

DropBox and *Google Drive* - both of which maintain multiple shadow copies of my work in the cloud.

Version Control

GitHub - which facilitated exploring different implementation strategies by forking my core code repository.

Chapter 3

Implementation

In this chapter I will give an overview of the implementation of the project, as well as investigating some of its more interesting features.

The chapter will be split into six parts:

Map generation

an explanation of the algorithm used to generate maps.

Simulation

an overview of the simulator, including graph generation and algorithm data.

Algorithms

each of the algorithms will be covered, with specific attention given to Block A*.

User Interface

an overview of the graphical user interface.

Testing

an overview of the methods used to test the code for correctness.

Data extraction

an explanation of how large volumes of data were extracted for statistical analysis.

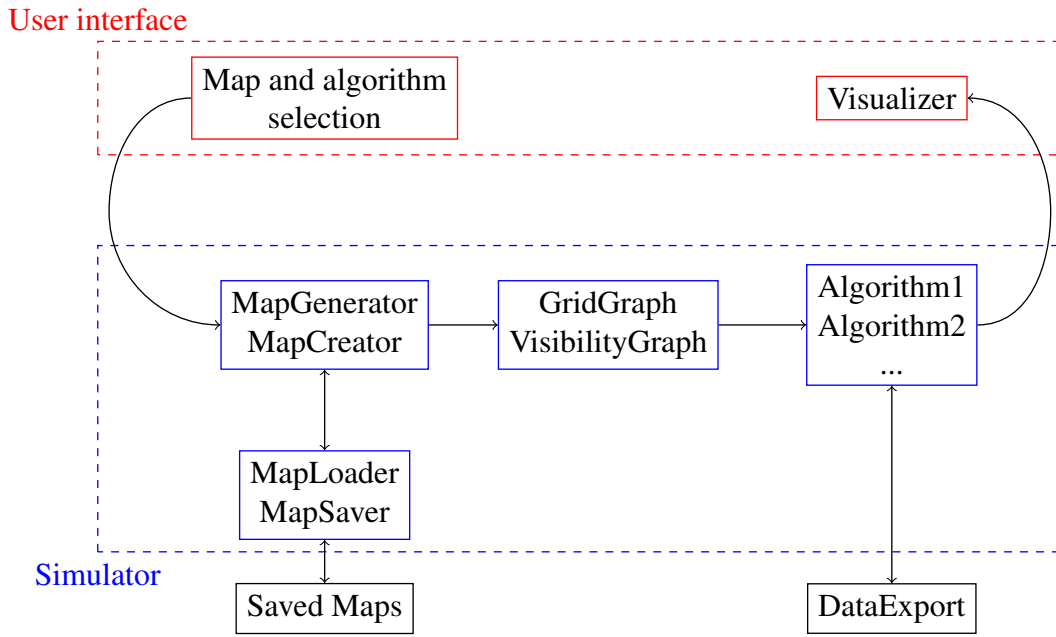


Figure 3.1: Flow of the user interface and simulator

3.1 Map generation

In order to reach reliable conclusions about the performance of the algorithms, I needed a way to generate large volumes of maps of a given resolution N , coverage percentage C and clustering D . I devised an algorithm to pseudo-randomly create such maps. The algorithm creates maps of varying clustering by loosely approximating the idea of potential fields.

The input to the algorithm includes an integer matrix $m_{i,j}$ of size N^2 - where each element has been initialised to value 1. At any point in the algorithm, $m_{i,j}$'s value represents its 'potential' - i.e. the chance that $m_{i,j}$ will be the next element set to 0 (i.e. blocked). The algorithm performs $C \times N^2$ iterations. In each iteration, it chooses a random number r between 0 and $\sum_{i,j} m_{i,j}$. It then traverses the matrix row-by-row until the sum of the elements it has seen is at least r . The potential of the element that has been reached is set to 0, and the potential of the surrounding elements is increased in a crude approximation of a potential field. The output is an integer matrix $m_{i,j}$ of size N^2 - where $C\%$ of the elements have value 0, which denotes a blocked cell, and the rest of the elements have value >0 , which denotes a free cell. The array can then be parsed into a Map by the Map constructor.

1	1	1	1
1	1	1	1
1	1	1	1
1	1	1	1

(a) Initialisation

1	1	1	1
3	5	3	1
5	0	5	1
3	5	3	1

(b) Iteration 1: $r = 5$

1	1	1	1
3	5	3	1
5	0	9	3
3	9	0	5

(c) Iteration 2: $r = 2$ Figure 3.2: Two iterations of GenerateMap with $R=4$ and $D=2$

Algorithm 6: GENERATEMAP

```

def GenerateMap( $m, C, D$ )
1  repeat
2     $r \leftarrow \text{random}(0, \sum_{i,j} m_{i,j})$ 
3     $i, j \leftarrow 0$ 
4    while  $r \geq 0$  do
5       $r \leftarrow r - m_{i,j}$ 
6      if  $i < R - 1$  then
7         $i \leftarrow i + 1$ 
8      else
9         $i \leftarrow 0$ 
10        $j \leftarrow j + 1$ 
11     SetAsBlocked( $i, j$ );
    until  $(C \times N^2)$  times
def SetAsBlocked( $m_{i,j}$ )
12   $m_{i,j} \leftarrow 0$ 
13  foreach  $m_{k,l}$  in horizontalOrVerticalNeighbour( $m_{i,j}$ ) do
14    if  $m_{k,l} \neq 0$  then
15       $m_{k,l} \leftarrow m_{k,l} + D$ 
16  foreach  $m_{k,l}$  in diagonalNeighbour( $m_{i,j}$ ) do
17    if  $m_{k,l} \neq 0$  then
       $m_{k,l} \leftarrow a_{k,l} + 2 \times D$ 

```

3.2 Simulation

3.2.1 Graph Generation

To create a graph from a map, the map was iterated over twice. The first iteration created nodes for all the coordinates that were valid node locations. The second iteration set the neighbours of each node. Checks for both of the iterations were performed by checking whether neighbouring blocks were blocked. See Figure 3.3.

As discussed in the Preparation chapter, along with grid-based graphs, I am also going to investigate the performance benefits of running path-finding algorithms on visibility graphs. The creation of visibility graphs depends on having a LineOfSight function.

Line of Sight

The line of sight algorithm is based on the pseudocode in [reference Theta* paper], which itself is a derivative of Bresenham's line drawing algorithm - though instead of choosing pixels (i.e. cells) to draw, it chooses cells to check whether they are blocked. Bresenham's algorithm is a useful framework as it avoids any floating point calculations when the start and end points are integers - this has dual benefits:

- the algorithm is fast
- the algorithm doesn't suffer from rounding errors inherent in floating-point calculations.

A notable alteration to the basic Bresenham algorithm is that instead of checking one cell per column (or one per row), the line of sight algorithm will check any cell that the real line passes through. This


```

1  ...
2  Coordinate[] diagonalRelativeCellCoordinates =
3      {new Coordinate(-1,-1), new Coordinate(0,-1),
4        new Coordinate(0,0), new Coordinate(-1,0)};
5
6  Node[][] graphArray2D =
7      new Node[map.getWidth()+1][map.getHeight()+1];
8  for(int j=0; j < map.getHeight()+1; j++) {
9      for(int i=0; i < map.getWidth()+1; i++) {
10         boolean isUnblockedAdjacentCell = false;
11         for(Coordinate c: diagonalRelativeCellCoordinates) {
12             try {
13                 if(!map.getCell(i+c.getX(),j+c.getY()).isBlocked()) {
14                     isUnblockedAdjacentCell=true;
15                 }
16             } catch (ArrayIndexOutOfBoundsException e) {}
17         }
18         if(!isUnblockedAdjacentCell) {
19             graphArray2D[i][j] = null;
20         } else {
21             graphArray2D[i][j] = new Node(new Coordinate(i,j));
22         }
23     }
24 }
25 ...

```

Figure 3.3: Code snippet showing Node creation for grid-based graphs

only requires a minor alteration.

For the purposes of clarity, the pseudocode presented in LINEOFSIGHT assumes the line of sight is in octant 1 - i.e. whose angle with the x-axis is between 0° and 45° . The full pseudocode can be found in Appendix ?. On this assumption, the key variables are:

x and y - integers that represent the coordinate of the cell being considered, which is always a cell that the line passes through.

f - a value that represents at what point the line intersects $x + 1$ with respect to the current y value. See Figure 3.4.

The algorithm starts at n_{start} . The current f -value is increased by dy every time x is increased, but decreased by dx if y is increased. If at any point $f = 0$ then the line intersects the bottom right-corner of the cell currently being considered, or if $y = dy$ then it intersects at the top left-corner. Therefore, if $f > 0$ then we need to check the $cell_{x,y}$ to see if it's blocked and if $f > dy$ then we additionally need to check $cell_{x,y+1}$.

Note: to disallow a line of sight through a 'diagonal blockage' (as introduced in section 2.1.3, a new **if** clause would be added after line 18 of LINEOFSIGHT:

if $f = 0 \wedge cell_{x,y}.isBlocked() \wedge cell_{x+1,y-1}.isBlocked()$ **then return false**.

To avoid the possibility of the final clause in the condition throwing an `ArrayIndexOutOfBoundsException`, an extra x -coordinate check or a try-catch block would also be required.

3.2.2 Algorithm Data

Each `MapInstance` has a `Map`, a `Graph` and an `AlgorithmData` object for each of the algorithms that have been run on that map.

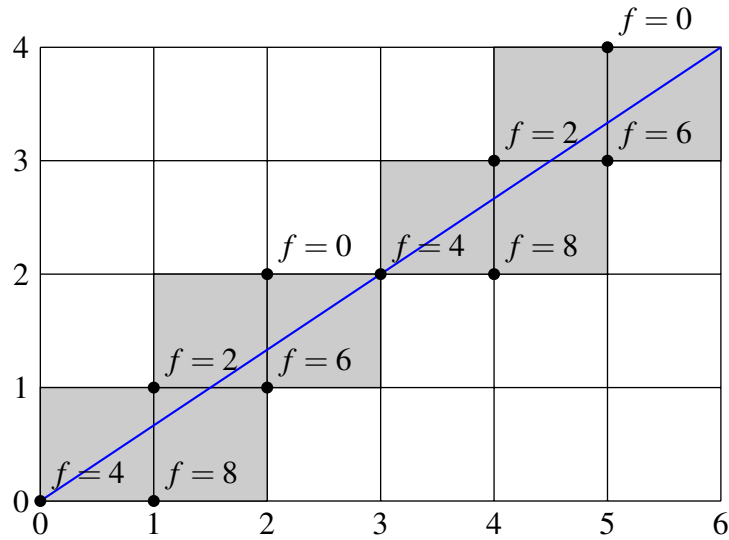


Figure 3.4: Line of sight algorithm

AlgorithmData is an abstract class to facilitate adding new algorithms in a modular way. It has a public method `go()`, as well as getter methods for all statistical data about each algorithm. `go()` calls the `getPath(n_{start}, n_{goal})` method, which returns n_{goal} if a path is found, or *null* otherwise. If *null* is returned, this is recorded in the statistics. If n_{goal} is returned, `go()` calculates some of the statistics as follows:

Path length

the sum of the Euclidean distances between each pair of consecutive nodes in the path.

Cumulative path angle

the sum of the scalar product between each pair of adjacent path segments in the path.

Graph calculation time

`System.nanoTime()` is used to calculate the duration between the call to `generateGraph()` and it returning.

Path calculation time

`System.nanoTime()` is used to calculate the duration between the call to `getPath()` and it returning.

Number of nodes expanded

a counter was incremented each time a node was expanded. It should be noted that in Block A*, node expansion is not the same as block expansion.

3.3 Algorithms

To emphasise the close relationships between the algorithms and allow for maximum code re-use, I organised the concrete instances of AlgorithmData in a hierarchy. This also ensured that any performance differences between algorithms was due to the different nature of each algorithm and not different implementation of similar concepts. See Figure 3.5.

Algorithm 7: LINEOF SIGHT

```

def LineOfSight ( $n_{source}, n_{goal}$ )
1    $x \leftarrow n_{source}.x$ 
2    $y \leftarrow n_{source}.y$ 
3    $x_{goal} \leftarrow n_{goal}.x$ 
4    $y_{goal} \leftarrow n_{goal}.y$ 
5    $dx \leftarrow x_{goal} - x$ 
6    $dy \leftarrow y_{goal} - y$ 
7    $f \leftarrow 0$ 
8   while  $x \neq x_{goal}$  do
9        $f \leftarrow f + dy$ 
10      if  $x \geq dx$  then
11          if  $cell_{x,y}.isBlocked()$  then
12              return false
13           $y \leftarrow y + 1$ 
14           $f \leftarrow f - 1$ 
15      if  $f \neq 0 \wedge cell_{x,y}.isBlocked()$  then
16          return false
17      if  $dy = 0 \wedge cell_{x,y}.isBlocked() \wedge cell_{x,y-1}.isBlocked()$  then
18          return false
19       $x \leftarrow x + 1$ 
20  return true

```

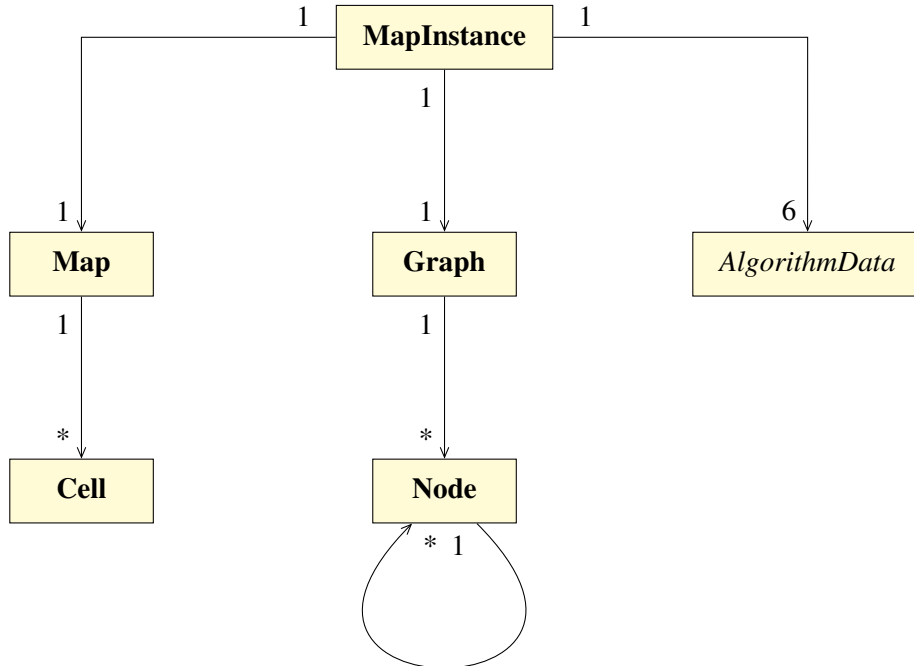


Figure 3.5: Composition of MapInstance

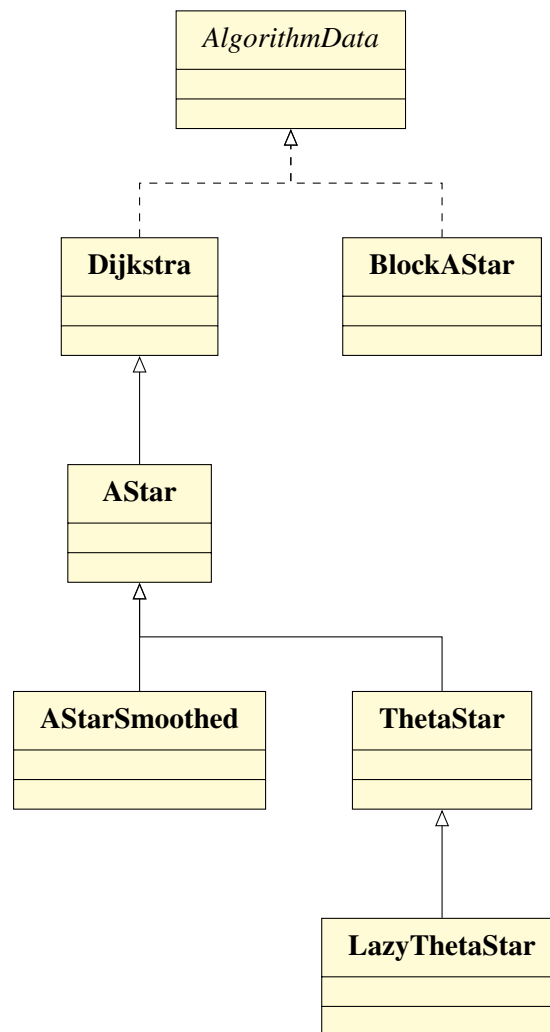


Figure 3.6: Inheritance structure of algorithms

3.3.1 Dijkstra's shortest paths

The implementation was based on the pseudo-code seen in section 2.2.1. Details of note include:

Open Set

the standard `java.util.PriorityQueue` would occasionally not return the node with the smallest g-value when `openSet.pop()` was called. This is a documented bug for large queues that occurs when a queue item has its priority¹ altered while residing in the `openSet`, so I had to manually `pop()`, `update()` and `re-add()` any node that needed to be updated when already in the `openSet`.

Closed set

the only two operations on the `closedSet` are adding and checking for membership. Therefore, a `HashSet` was used for its average case $O(1)$ insertion and search speed.

Extensibility

to allow a clear algorithm hierarchy, I needed to add a call to `initialise(n)` whenever a node n was popped from the `openSet`, and a `postProcessing(n)` step before the node is returned. In Dijkstra these have empty method bodies, but some of the algorithms that inherit from Dijkstra will override these methods.

3.3.2 A*

The implementation was based on the pseudo-code seen in section 2.2.2: A* inherits from Dijkstra, and overrides the `updateCost()` function.

3.3.3 A* with post-smoothing

The implementation was based on the pseudo-code seen in section 2.2.3: A* inherits from A* with post-smoothing and implements its post-smoothing by overriding the `postProcessing()` function.

3.3.4 Basic θ^* and Lazy θ^*

The implementation was based on the pseudo-code seen in section 2.2.4 and 2.2.5:

Basic θ^*

inherits from A*, and overrides the `updateCost()` function

Lazy θ^*

inherits from Basic θ^* , and overrides the `initialiseNode()` and `updateCost()` functions

3.3.5 Block A*

Local Distance Database

The first challenge was to obtain an LDDb. Since there was no publicly available library containing the LDDb, I had to create my own. Although it would have been possible to create the entries manually for the LDDb for block sizes of 2×2 cells, this would certainly not have been feasible for block sizes that were any larger, since for blocks of size $n \times n$, there will be 2^{n^2} possible blocks, each with $4.n.4.(n-1)$ pairs of ingress and egress coordinates, which gives over 12 million calculations for a

¹In this case, the Node's g-value

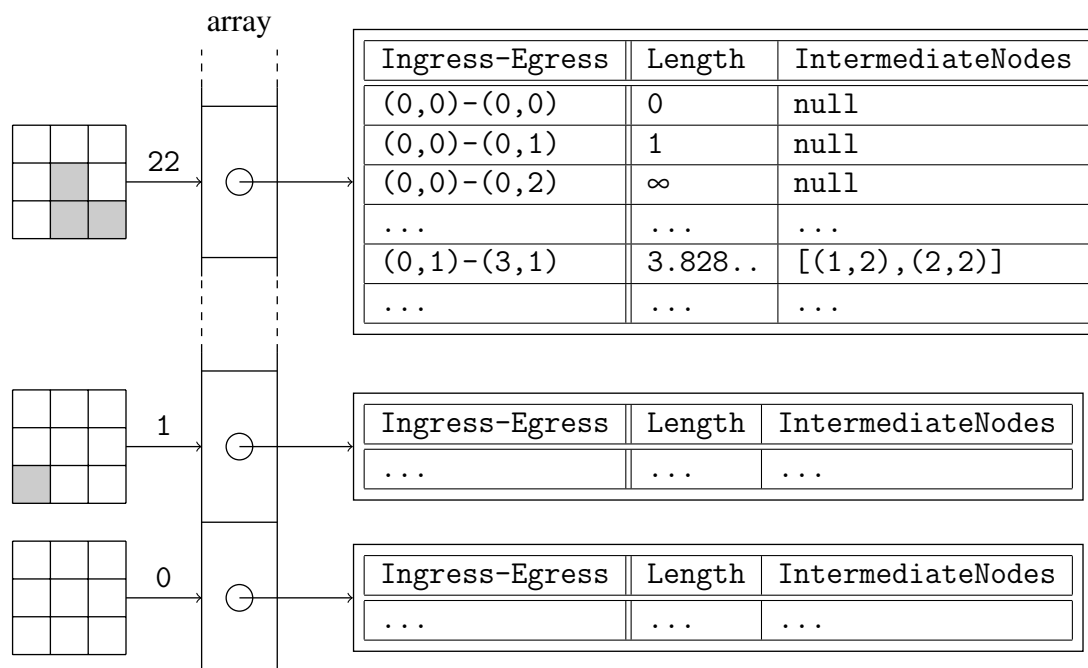


Figure 3.7: Extract of the LDDb for block size of 3×3 - array of `HashTable<PairOfCoords, Pair<double, List<Coord>>>`

block size of 4×4 .

I calculated the shortest paths using A* over visibility graphs, which gives provably optimal paths.

For a given block size, I required:

for every possible block of size $n \times n$

for every possible ingress-egress coordinate pair

(a) the shortest path between that pair

(b) a list containing every intermediate coordinate on that path.

It was necessary to put these entries into a database structure that would be compact in memory and fast to load and query. Since these constraints were fundamental to the operation of the algorithm, any library or 3rd party database implementation could not be guaranteed to be sufficiently specialised for the task, so I implemented my own database using arrays and hash tables to ensure optimal performance and minimum space wastage.

All queries to the database would need to identify a block and pair of coordinates as an input, and would receive either the length of the shortest path between those nodes, or a list of the intermediate nodes on the shortest path between those two nodes. I devised a simple bitwise encoding scheme to represent blocks as integers that would be much more space and time efficient than storing Map objects in the database to use as comparison: each cell in the block is represented by a bit in the integer, and that bit is set if the corresponding cell is blocked. The 32 bits of the integer are sufficient for all the block sizes that are worth considering: 2×2 to 4×4 .¹

Using this scheme, the underlying implementation of my database was an array of HashMaps - one HashMap per block, with the array indexed by the code of the block. The HashMap mapped a key: a Pair of ingress-egress Coordinates, to a value: a Pair consisting of the length (as a double) of the shortest path and an ArrayList of the Coordinates on the path.

```

1 PairOfCoords(Coordinate c1, Coordinate c2, int blockSize) {
2     this.blockSize = blockSize;
3     this.coordCode = (byte)
4         (c1.getX() +
5          c1.getY() * (blockSize+1) +
6          c2.getX() * (blockSize+1)*(blockSize+1) +
7          c2.getY() * (blockSize+1)*(blockSize+1)*(blockSize+1));
8 }

```

Figure 3.8: Encoding scheme for PairOfCoords in the compressed LDDB

```

1 int getListCode(List<Coordinate> intermediateNodes) {
2     int listCode = 0;
3     for(Coordinate c : intermediateNodes) {
4         listCode = listCode | c.getX();
5         listCode = listCode << 3;
6         listCode = listCode | c.getY();
7         listCode = listCode << 3;
8     }
9     listCode = listCode >>> 3; //undoes line 7 on final loop
10    return listCode;
11 }

```

Figure 3.9: Encoding scheme for List<Coordinate> in the compressed LDDB

This implementation was sufficient but unsatisfactory, as the database sizes were unnecessarily large. This would cause slow loading to memory, less of the LDDB stored in cache and more likelihood of thrashing. Therefore, I compressed the database using bitwise encoding schemes:

- each Pair of ingress-egress Coordinates can be represented with a unique integer representable in a byte's worth of space - so a coding scheme was devised to create an efficient hash function for the PairOfCoordinate class.
- the List of intermediate Coordinates can be represented by a code that fits into the 32 bits of an integer: since the maximum number of intermediate nodes on a shortest path in a sub map of size up to 4×4 is four², and the range of x and y in the Coordinate is 0-4, therefore 6 bits can be used for each Coordinate - 3 for x , 3 for y , which is a maximum total of 24 bits.

These compression techniques allowed me to reduce the size of the databases by about 50%:

Submap size	Size of uncompressed LDDB	Size of compressed LDDB ^{3,4}
2×2	33KB	17KB
3×3	2.2MB	1.1MB
4×4	485MB	192MB

I will investigate the performance of the uncompressed and compressed LDDBs in the Evaluation section.

¹The exponential increase in size of LDDB would mean that the LDDB for block sizes of 5×5 or larger take so long to search that performance benefit diminishes.

²Obtained by experimental results.

³These savings could be improved further for sub maps of size 2×2 and 3×3 by using specific data-types depending on the submap size, but this was unnecessary, as the LDDB size was very small compared to the total available memory for sub maps of size 2×2 and 3×3 , and the LDDB only needs to be loaded into memory once per run of the program.

⁴These sizes were significantly larger than those found in Yapp's paper, but his agents were only allowed to travel

Special case blocks: $block_{start}$ and $block_{goal}$

$block_{start}$ and $block_{goal}$ are treated as special cases by Block A*. On initialisation:

- the g -values of the boundary nodes of $block_{start}$ are set by manually computing shortest paths from n_{start} to each boundary node, instead of looking them up in the LDDb as with other blocks. This is because n_{start} is not guaranteed to be on the boundary of a block, and the LDDb only holds details for routes between boundary nodes
- the h -values of the boundary nodes of $block_{goal}$ are set manually for similar reasons.
- a check is made as to whether n_{start} and n_{goal} are in the same block - i.e. $block_{start}$ equals $block_{goal}$. If so, the shortest path between the two is computed manually.

To compute these values, I created a visibility graph for the relevant block and computed the shortest path to each boundary node. Having run some preliminary tests, it became clear that on small maps these initialisations took up to 50% of the run-time of the algorithm [insert some statistics of test runs]. Therefore, I decided to implement two alternative, extended forms of the LDDb:

Semi-extended contains details of shortest paths from *any* node to *any* boundary node in a block.

Fully-extended contains details of shortest paths from *any* node to *any* node in a block.

The increase in the size of the LDDb in comparison to the original is:

Submap size	Semi-extended	Fully-extended
2×2	12.5%	26.5%
3×3	33.3%	77.8%
4×4	56.3%	144.1%

The use of the fully-extended LDDb allows the initialisation of $block_{start}$ and $block_{goal}$ and the $block_{start}$ equals $block_{goal}$ case to use the LDDb, whereas using semi-extended LDDb requires that the $block_{start}$ equals $block_{goal}$ case finds the shortest route manually.¹ The performance benefits of these different forms of the LDDb will be investigated in the Evaluation chapter.

Block A* - algorithm

The core of the algorithm was a straightforward implementation, once the details of it were understood. The most challenging part of the implementation was the traceback stage, which was unspecified in the paper. The traceback stage extracts the coordinates of the path from the Graph structure. This requires finding and inserting into a list:

- the optimal path from n_{goal} to the optimal ingress node of $block_{goal}$
- the nodes on the optimal path that share block boundaries
- the intermediate nodes where the optimal path traverses a block
- the optimal route from the optimal egress node of $block_{start}$ to n_{start}

whilst omitting from the list:

- all but one of any nodes in this path that lie in the same physical location but belong to different blocks: this occurs when the path crosses block boundaries.

horizontally and vertically, so only 4 bits were needed to store the path length value.

¹Note that when using the semi-extended LDDb, $block_{goal}$ needs to reverse the ingress-egress coordinates before querying the LDDb.

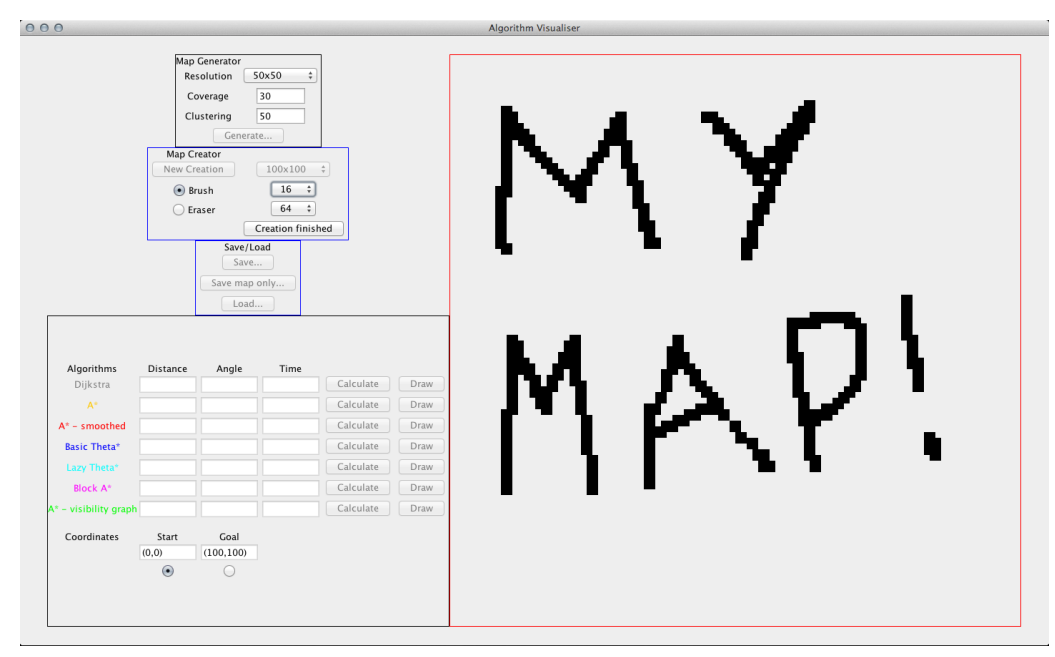


Figure 3.10: User interface in map creation mode

3.4 User Interface

The user interface was built using Java's Swing toolkit. The UI allows the user to:

Generate a map of a given resolution, coverage percentage and clustering.

Create a map of a given resolution using an interactive editor that has brushes and erasers of differing sizes. If the UI is in map creation mode then UI uses `MouseListeners` to detect when the mouse is being dragged over the map. The array that represents the map is updated according to the size of the brush and the resolution of the map, and then the `repaint()` method of the `JPanel` is called. Once the creation is complete, the array is passed as a parameter to a `Map` constructor, and the `Map` object is passed to the `Simulator`.

Save a map with or without paths and statistics.

Load a map with or without paths and statistics.

Choose start and end points for the paths on the current map. If the UI is not in map creation mode then `MouseListeners` are used to detect where on the screen the mouse is clicked. The start and end points can be set anywhere on the map.

Calculate paths for each of the 6 algorithms on the current map. If there is no path then `NoPath` will be displayed.

Display paths once the path has been calculated

3.5 Testing

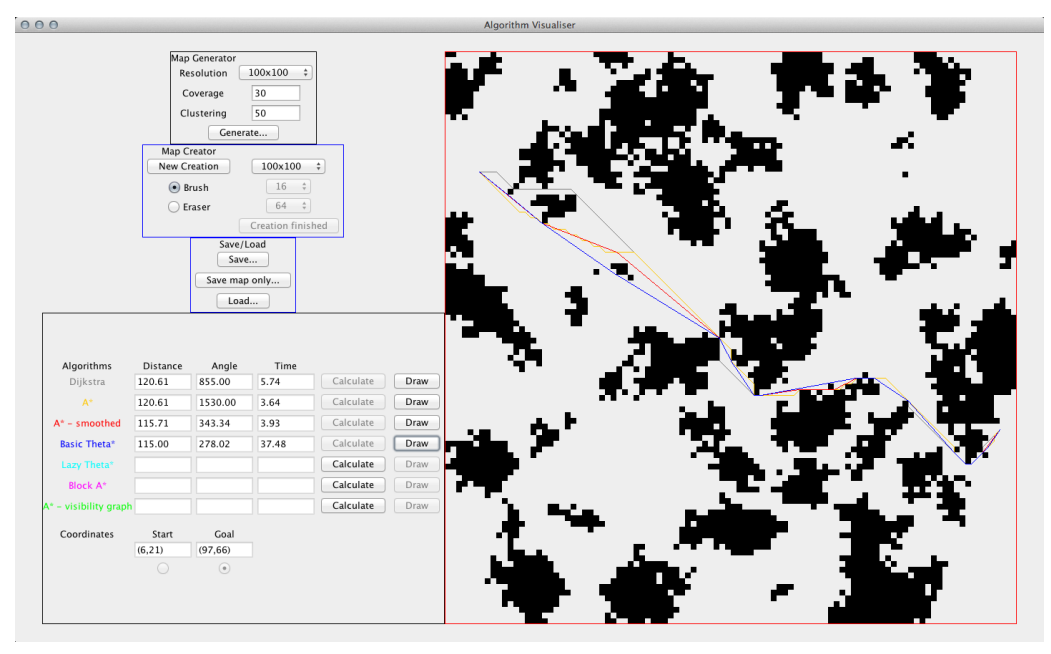


Figure 3.11: User interface displaying paths for a generated map

Map	Algorithm	PathTime	TotalLength	TotalAngle
Map 1	Dijkstra	852.926976	160.752308679	495.0000736525
Map 1	AStar	169.831936	160.752308679	855.0000688228
Map 1	AStarVisibility	5.334016	151.6768359881	102.262577962
Map 1	AStarSmoothed	85.908992	154.9127142045	165.3889464848
Map 1	ThetaStar	30.230016	151.7637935619	122.1222152092
Map 1	LazyThetaStar	32.45824	151.7637935619	122.1222152092
Map 1	BlockAStar	10.85792	152.8713149055	566.7296353554
Map 2	Dijkstra	615.220992	161.3380951166	810.0000676154

Figure 3.12: Extract from a CSV file exported by DataExtract

3.6 Data extraction

To harvest enough data to do meaningful statistical analysis of the performance of the algorithms, the simulator was designed to have a simple API that allowed both a UI to be bolted on and scripts that could bypass the UI altogether. I used the open source package `CSVWriter` to write the data obtained to CSV files, as these are the standard input for *R* based statistical analysis.

```
1 void generateMaps
2 (int size, int coverage, int clustering, int numberOfMaps) {
3
4     for(int i=0;i<numberOfMaps;i++) {
5         saveMapOnly(size+"/"+coverage+"/"+clustering+"/"+i+".ser",
6             new MapInstance(
7                 MapGenerator.generateMap(size,size,coverage,clustering)));
8     }
9 }
```

Figure 3.13: Code snippet from DataExtraction showing automation of Graph creation

Chapter 4

Evaluation

Chapter 5

Conclusion