

1           Comparison of Semi-Implicit and Fully-Implicit  
2           Methods for a Highly Degenerate Diffusion-Reaction  
3           Equation Coupled with an ODE

4           by

5           Eric M. Jalbert

6           A Thesis  
7           presented to  
8           The University of Guelph

9           In partial fulfilment of requirements  
10          for the degree of  
11          Master of Science  
12          in  
13          Applied Mathematics

14           Guelph, Ontario, Canada

15           © E.M. Jalbert, January, 2015

# <sup>16</sup> **Contents**

<sup>17</sup> <b>1</b>	<b>Introductions</b>	<b>1</b>
<sup>18</sup> 1.1	Background . . . . .	1
<sup>19</sup> 1.2	Objectives . . . . .	1
<sup>20</sup> 1.3	Outline . . . . .	2
<sup>21</sup> <b>2</b>	<b>Model Definition</b>	<b>3</b>
<sup>22</sup> 2.1	Model Description . . . . .	3
<sup>23</sup> 2.2	Nondimensionalization . . . . .	6
<sup>24</sup> 2.3	Parameters . . . . .	8
<sup>25</sup> <b>3</b>	<b>Numerical Methods</b>	<b>9</b>
<sup>26</sup> 3.1	Discretization . . . . .	9
<sup>27</sup> 3.2	Solving Technique . . . . .	10
<sup>28</sup> 3.3	Computational Setup . . . . .	17
<sup>29</sup> 3.4	Method Validation . . . . .	17
<sup>30</sup> 3.5	Comparison of Semi-implicit and Fully-implicit Method . . . . .	24
<sup>31</sup> <b>4</b>	<b>Simulation Results</b>	<b>27</b>
<sup>32</sup>	<b>Complete Bibliography</b>	<b>28</b>

<sup>33</sup> **Chapter 1**

<sup>34</sup> **Introductions**

<sup>35</sup> **1.1 Background**

- <sup>36</sup> • C. thermocellum biology and the connection to ethanol production. Look at USA data and the fact that C.thermo is a cellulolytic bacteria.
- <sup>38</sup> • Talk about existing models for biofilms, such as: suspended cultures, reactor scale-models, and Cellular Automata.
- <sup>40</sup> • Look at PDE with Volume filling and compare the model we will use here to the traditional biofilm model.
- <sup>42</sup> • Describe the Numerics of PDE's, such as semi-implicit methods and the different work that has been done to solve these kinds of problems. Mention the problems those methods would have with our system and mention where this work fits in.

<sup>45</sup> **1.2 Objectives**

<sup>46</sup> This will be one to two paragraphs for each objective.

- <sup>47</sup> 1. Numerical development, implementation and validation of a fully implicit method

- 48        2. The comparison fo the fully implciit method with the semi-implicit methods
- 49        3. Simulation work; see if we can better understand the mechanisms of the bacteria by extracting
- 50           some interesting observations from our simulations.

51        **1.3   Outline**

- 52        This will be a short section that describes what each section after this one will be covering and how it
- 53        helps to accomplish the objectives, either in bullet point or one sentence for each section.
- 54        When you write this section, also make sure that you describe how each of these sections/chapters
- 55        relates to the objectives that you formulated, and how these sections relate to each other

56 **Chapter 2**

57 **Model Definition**

58 **2.1 Model Description**

59 The model used for simulations is based on the deterministic biofilm model developed in Eberl et al.  
60 (2001), which was designed for modelling the development of spatially heterogenous biofilm struc-  
61 tures. Since *C.Thermocellum* grows as a monolayered biofilm and consumes a solid carbon fiberous  
62 substrate, there are mechanical differences between our system. For our model there are a number of  
63 assumptions that it must satisfy:

- 64 1. The growth of sessile biomass is inhibited by the nutritional value of the substrate and by the  
65 availability of colonizable space.
- 66 2. The number of cells on a square unit of substratum is limited to a finite value because *C.Thermocellum*  
67 forms only a thin monolayer.
- 68 3. Biomass remains sessile unless the density approaches close to the physical limit. Near the  
69 physical limit, biomass must diffuse immediatly. Biomass density should never reach the phys-  
70 ical limit.
- 71 4. The carbon fiberous substrate consumed as nutrition for biomass growth is the substratum to  
72 which the biofilm attachs. However, fully consumed substrate does not remove the substratum.

73 Instead the biofilm propagation is forced by the lack of nutrients as the substrate concentration  
 74 decreases. This also implies that the substrate does not diffuse.

75 5. Bacteria production is primarily limited by available substrate. With an over abundance of  
 76 substrate, the growth of bacteria would quickly consume the excess and force the limiting factor  
 77 of growth to be substrate availability.

78 6. A fraction of cells dies constantly.

79 7. Substrate consumption is exactly related to biomass growth .

80 Assumption 1, 2, 3, 5, 6, and 7 are similar to those made in Eberl et al. (2001). The main difference  
 81 here is 4; our substrate is sessile. With a sessile substrate, there is no diffusion for the substrate  
 82 concentration. Another difference is that there is no fluid mechanics here. Instead our biofilm rests  
 83 on a solid physical substratum.

84 Consider a spatial region,  $\Omega$ , here the model formed must relate the variables:

$t \geq 0$  time as an independent variable

$x \in \Omega$  spatial coordinate as an independent variable

$M(t, x)$  biomass density fraction as a dependent variable

85  $C(t, x)$  substrate density as a dependent variable

$d(M)$  diffusion coefficient of as a variable model parameter

$f(C)$  biomass production rate as a variable model parameter

$g(C)$  substrate consumption rate as a variable model parameter

86 The biomass density fraction represents the current density of biomass divided by the maximum  
 87 biomass density,  $M_\infty$ .

88 From the above assumptions, a PDE-ODE-coupled system that models *C. Thermocellum* growth can  
 89 be purposed as,

$$M_t = \nabla_x (d(M) \nabla_x M) + f(C)M \quad (2.1)$$

$$C_t = -g(C)MM_\infty \quad (2.2)$$

90 where

$$91 \quad d(M) = d \frac{M^\alpha}{(1 - M)^\beta} \quad (2.3)$$

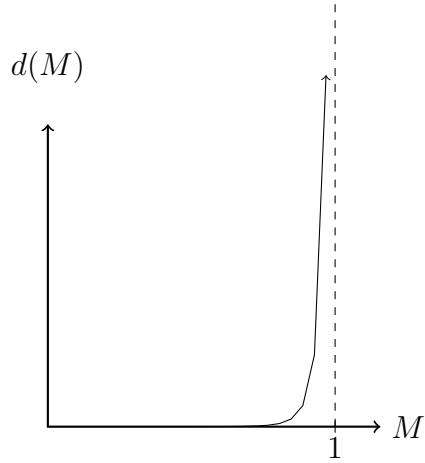
$$92 \quad 93 \quad f(C) = u \frac{C}{k + C} - n \quad (2.4)$$

$$94 \quad 95 \quad g(C) = y \frac{C}{k + C} \quad (2.5)$$

96 Here we have a pair of equations, (2.1) and (2.2), that represent the biomass density and substrate  
 97 concentration respectivly. The terms in (2.1) correspond to the density-dependent diffusion term,  
 98  $\nabla_x (d(M) \nabla_x M)$ , and a biomass production term,  $f(C)M$ . This satisfies assumption 1 since the only  
 99 factors effecting the biomass density is growth from nutrient converton and diffusion from spatially-  
 100 full colonized space. Coinciding with assumption 3 the density-dependent diffusion equation (2.3)  
 101 has a near-zero value until when  $M \rightarrow 1$ , which leads to  $d(M) \rightarrow \infty$ ; this can be seen in Figure 2.1.  
 102 The production rate is a simple Monod kinetic growth term with a constant death rate term to match  
 103 with assumption 5 and 6. Monod kinetic growth was selected since it matchs the growth of baterial  
 104 when limited by availble nutrients.

105 The substrate concentration equation (2.2) is formed from the combination of a consumption function  
 106 (2.5), the current biomass density, and the maximum density of biomass. Agreeing with assumption  
 107 4 there is no diffusion term, only the bacteria consuming the substrate nutrients will lower the con-  
 108 centration. Consumption function (2.5) is based on the growth of the biomass by a scalar multiplier  
 109 and thus satisfying assumption 7. The,  $M_\infty$ , relates the density fraction  $M$  to the physical density  $C$ .

110 The dimensions of the parameters and variables are in Tabel 2.1.



**Figure 2.1:** A graph of  $d(M) = d \frac{M^\alpha}{(1-M)^\beta}$  showing the way diffusion increases asymptotically as  $M \rightarrow 1$ .

Variable/Parameter	Dimensions
$t$	[days]
$x$	[meters]
$M$	[−]
$C$	[ $\frac{\text{grams}}{\text{meters}^3}$ ]
$d$	[ $\frac{\text{meters}^2}{\text{days}}$ ]
$\alpha$	[−]
$\beta$	[−]
$u$	[ $\text{days}^{-1}$ ]
$k$	[ $\frac{\text{grams}}{\text{meters}^3}$ ]
$y$	[ $\text{days}^{-1}$ ]
$n$	[ $\frac{\text{grams}}{\text{meters}^3 \cdot \text{days}}$ ]
$M_\infty$	[ $\frac{\text{grams}}{\text{meters}^3}$ ]

**Table 2.1:** List of parameters and their dimensions

## 111 2.2 Nondimensionalization

112 To help facilitate the analysis of this system, the full removal of all physical units is preferred and  
 113 so we nondimensionalize the parameters. Here the parameters used are: the biomass growth rate,  $u$ ;  
 114 the length of the region,  $L$ ; and the maximum density for biomass and substrate,  $M_\infty$  and  $C_\infty$ . From  
 115 using the following parameter changes, the system can be made unitless.

$$\chi = \frac{x}{L} \implies Ld\chi = dx \quad (2.6)$$

$$\tau = ut \implies \frac{1}{u}d\tau = dt \quad (2.7)$$

$$\mathcal{C} = \frac{C}{C_\infty} \quad (2.8)$$

$$\delta = \frac{1}{uL^2}d \quad (2.9)$$

$$\kappa = \frac{k}{C_\infty} \quad (2.10)$$

$$\nu = \frac{n}{uC_\infty} \quad (2.11)$$

$$\gamma = \frac{M_\infty}{C_\infty}y \quad (2.12)$$

<sup>116</sup> Using these, (2.1) and (2.2) can be simplified and nondimensionalized into,

$$M_\tau = \nabla_\chi (D(M)\nabla_\chi M) + F(\mathcal{C})M \quad (2.13)$$

$$\mathcal{C}_\tau = -G(\mathcal{C})M, \quad (2.14)$$

<sup>117</sup> where,

$$\begin{aligned} D(M) &= \delta \frac{M^\alpha}{(1-M)^\beta} \\ F(\mathcal{C}) &= \frac{\mathcal{C}}{\kappa + \mathcal{C}} - \nu \\ G(\mathcal{C}) &= \gamma \frac{\mathcal{C}}{\kappa + \mathcal{C}}. \end{aligned} \quad (2.15)$$

<sup>119</sup> with only  $\delta, \kappa, \nu, \gamma$  as model parameters. For convenience, we henceforth use

$$C := \mathcal{C}, \quad x := \chi, \quad t := \tau. \quad (2.16)$$

**121 2.3 Parameters**

122 Each of the dimensionless parameters in (2.15) have a biological representation based on the transfor-  
123 mations done. The parameter  $\delta$  is the dimensionaless constant for diffusion. It affects the change in  
124 biomass from adjacent biomass sources, a greater  $\delta$  results in a greater change. The parameter  $\kappa$  is the  
125 half-saturation point, it is exactly the value for which substrate concentration results in 0.5-optimum  
126 growth rate. Parameter  $\nu$  is the death rate of the biomass. Specifically, it is the ratio of biomass  
127 growth to death, representing the fraction of biomass density that perishes from natural causes or a  
128 lack of substrate. Lastly,  $\gamma$  is the yield ratio. It signifies the ratio of substrate consumed to biomass  
129 growth. Here, a larger  $\gamma$  value results in more substrate being consumed to produce the same amount  
130 of biomass.

131 With (2.13) being reduced to four parameters the numerical analysis become more simiplified while  
132 still retaining the same significance in results.

<sup>133</sup> **Chapter 3**

<sup>134</sup> **Numerical Methods**

<sup>135</sup> **3.1 Discretization**

<sup>136</sup> In order to find the solution for (2.13) spatial and temporal discretizations must be made. First the  
<sup>137</sup> equations are discretized in time,

$$\frac{M^{k+1} - M^k}{\Delta t} = \nabla_x(D(M^{k+1})\nabla_x M^{k+1}) + F(C^{k+1})M^{k+1}, \quad (3.1)$$

$$\frac{C^{k+1} - C^k}{\Delta t} = \frac{h}{2}(G(C^{k+1})M^{k+1} + G(C^k)M^k). \quad (3.2)$$

<sup>141</sup> Here, (3.1) follows the ideas of the Backwards Euler Method; (3.2) follows Trapezoidal Rule. The  
<sup>142</sup> index variable  $k$  has also been introduced in (3.1 - 3.2) such that  $M^k(x) \approx M(t^k, x)$ , allowing an  
<sup>143</sup> approximation at a certain time,  $t^k$ , to be used; this reduces the dimensionality of the problem.

<sup>144</sup> For this system, the region of consideration will be a rectangular region,  $\Omega$ . This region has Neumann  
<sup>145</sup> boundary conditions,  $\frac{\partial M}{\partial x} = \frac{\partial C}{\partial x} = 0, \forall x \in \partial\Omega$ . Now, only (3.1) requires spatial considerations since  
<sup>146</sup> the substrate does not diffuse across the region. The spatial discretization will be through the Finite  
<sup>147</sup> Difference Method as described in Saad (2003). Here, a uniform  $n \times m$  grid is used to discretize  $\Omega$ .  
<sup>148</sup> Since all the calculations will be done on the grid intersections the discretization will be grid-point  
<sup>149</sup> based. This means that a  $n \times m$  grid implies there are  $(n - 1) \times (m - 1)$  grid boxes. The distance

150 between grid points is the same in both  $x_1$  and  $x_2$  dimensions; we have  $\Delta x_1 = \Delta x_2$ . A five-point  
 151 stencil is used to approximate the solution of (3.1) at each grid point. To index the grid point,  $i$  and  
 152  $j$  are used such that  $M_{i,j}^k \approx M(t^k, x_{1,i}, x_{2,j})$ . To account for the dependency on neighbouring grid  
 153 points, we introduce  $\sigma$  as the index pair from the set

$$154 \quad \mathcal{N}_{ij} = \{n_{ij}, e_{ij}, s_{ij}, w_{ij}\}. \quad (3.3)$$

155 where,

$$156 \quad \begin{aligned} n_{ij} &= \begin{cases} (i, j+1) & \text{if } j < m \\ (i, j-1) & \text{if } j = m \end{cases} & e_{ij} &= \begin{cases} (i+1, j) & \text{if } i < n \\ (i-1, j) & \text{if } i = n \end{cases} \\ s_{ij} &= \begin{cases} (i, j-1) & \text{if } j > 0 \\ (i, j+1) & \text{if } j = 0 \end{cases} & w_{ij} &= \begin{cases} (i-1, j) & \text{if } i > 0 \\ (i+1, j) & \text{if } i = 0 \end{cases} \end{aligned} \quad (3.4)$$

157 With  $\mathcal{N}_{ij}$  and  $\sigma$  we can account for the difference in boundary points and interior points.

158 The equation for (3.1), after spatial discretization, is

$$159 \quad \frac{M_{i,j}^{k+1} - M_{i,j}^k}{\Delta t} = \frac{1}{\Delta x^2} \sum_{\sigma \in \mathcal{N}_{ij}} \left( \frac{D(M_{\sigma}^{k+1}) + D(M_{i,j}^{k+1})}{2} \right) \cdot (M_{\sigma}^{k+1} - M_{i,j}^{k+1}) + F(C_{i,j}^{k+1}) M_{i,j}^{k+1} \quad (3.5)$$

160 For (3.5), the arithmetic mean of the diffusion function,  $D$ , is taken because of the steep gradient  
 161 at the interface. The alternative would be to use  $D(M_{i+\frac{s}{2},j+\frac{r}{2}}^{k+1})$ , however in some cases we have  
 162  $M_{i+\frac{s}{2},j+\frac{r}{2}}^{k+1} = 0$  which would result in  $D(0) = 0$  and thus nullify the effect of the spatial diffusion.  
 163 Taking the arithmetic mean eliminates this result because the average value of  $D$  would not be zero  
 164 at the interface.

## 165 3.2 Solving Technique

166 Now there exist equations for which  $C$  and  $M$  can be solved, (3.2) and (3.5) respectively. Using  $C^k$  and  
 167  $M_{i,j}^k$  as approximations of the solutions for (2.13) will allow the system to be solved by computing

<sup>168</sup>  $C^{k+1}$  and  $M_{i,j}^{k+1}$ . However, there are complications with trying to get an explicit formula for  $M_{i,j}^{k+1}$   
<sup>169</sup> from (3.5) because of the dependency on  $M$  in  $D(M)$ . To remedy this, a fixed point iteration is  
<sup>170</sup> introduced. In a single time step, the solutions for  $M$  and  $C$  can be solved using the previous time  
<sup>171</sup> step solution in the follow manner:

$$\frac{M_{i,j}^{(p+1)} - M_{i,j}^k}{\Delta t} = \frac{1}{\Delta x^2} \sum_{(s,r) \in \mathbb{A}} \left( \frac{D(M_{i+s,j+r}^{(p)}) + D(M_{i,j}^{(p)})}{2} \cdot (M_{i+s,j+r}^{(p+1)} - M_{i,j}^{(p+1)}) \right) + F(C_{i,j}^{(p)}) M_{i,j}^{(p+1)} \quad (3.6)$$

$$\frac{C^{(p+1)} - C^k}{\Delta t} = \frac{-1}{2} (G(C^{(p+1)}) M^{(p+1)} + G(C^k) M^k) \quad (3.7)$$

<sup>175</sup> where  $(p) \in (0, 1, \dots, P)$ . Note, that the equation for  $M_{i,j}^{(p+1)}$  shown in (3.6) refers to the interior  
<sup>176</sup> points only. A similar change is done for the boundary points but is not shown due to its complexity.  
<sup>177</sup> It is important to show explicitly that the purpose of the fixed point iteration is to link two distinct  
<sup>178</sup> times with  $P$  solutions in between them, such that:

$$\begin{aligned} M^{(p=0)} &= M^k, & M^{(p=P)} &= M^{k+1}, \\ C^{(p=0)} &= C^k, & C^{(p=P)} &= C^{k+1}. \end{aligned} \quad (3.8)$$

<sup>180</sup> In this fixed point format, given by (3.6 - 3.7), the equations can be rearrange and solved by conven-  
<sup>181</sup> tional methods.

<sup>182</sup> For (3.6), a linear system of equations can be created following Saad (2003). For each grid point  $(i, j)$   
<sup>183</sup> a linear system exists, defined as:

$$\begin{aligned} \frac{M_{i,j}^k}{\Delta t} &= \sum_{(i,j) \in \mathbb{A}} \left( \frac{D(M_{i+s,j+r}^{(p+1)}) + D(M_{i,j}^{(p+1)})}{2\Delta x^2} \cdot M_{i+s,j+r}^{(p+1)} \right) \\ &+ \left( \sum_{(i,j) \in \mathbb{A}} \left( \frac{D(M_{i+s,j+r}^{(p+1)}) + D(M_{i,j}^{(p+1)})}{2\Delta x^2} \right) - F(C_{i,j}^{(p)}) + \frac{1}{\Delta t} \right) M_{i,j}^{(p+1)}. \end{aligned} \quad (3.9)$$

<sup>185</sup> From (3.9), a five-diagonal matrix can be created defined as,

$$\text{186 } A = \begin{pmatrix} a_{i,j} & a_{i+1,j} & & a_{i,j+1} & & \\ a_{i-1,j} & \ddots & \ddots & & \ddots & \\ & \ddots & \ddots & \ddots & & \ddots \\ a_{i,j-1} & & a_{i-1,j} & a_{i,j} & a_{i+1,j} & a_{i,j+1} \\ & \ddots & & \ddots & \ddots & \ddots \\ & & a_{i,j-1} & a_{i-1,j} & a_{i,j} & a_{i+1,j} & a_{i,j+1} \\ & & & \ddots & \ddots & \ddots & \ddots \\ & & & & \ddots & \ddots & \ddots & a_{i+1,j} \\ & & & & a_{i,j-1} & a_{i-1,j} & a_{i,j} & \\ & & & & & \ddots & \ddots & \ddots \\ & & & & & & \ddots & \ddots & a_{i+1,j} \\ & & & & & & a_{i,j-1} & a_{i-1,j} & a_{i,j} \end{pmatrix} \quad (3.10)$$

<sup>187</sup> where each  $a_{i,j}$  is the coefficient based on (3.9).

<sup>188</sup> Solving (3.10) can be done by use of the Conjugate Gradient method provided that certain conditions  
<sup>189</sup> are satisfied.

<sup>190</sup> **Proposition 3.2.1.** *The matrix  $A$ , (3.10), is positive definite and symmetric when  $\frac{1}{F(C_{i,j}^{(p)})} < \Delta t$ .*

<sup>191</sup> *Proof.* Matrix  $A$  is positive definite if all the eigenvalues are positive. Using the Circle theorem  
<sup>192</sup> described by Geršgorin (1931), the eigenvalues can be shown to be positive if, independently on all  
<sup>193</sup> rows, the sum of the off-diagonals values is less than the diagonal value. This can be verified. From  
<sup>194</sup> (3.9) it can be said that,

$$\text{195 } \sum_{(i,j) \in \mathbb{A}} \left( \frac{D(M_{i+\frac{s}{2}, j+\frac{r}{2}}^{(p)})}{\Delta x^2} \right) < \left( \sum_{(i,j) \in \mathbb{A}} \left( \frac{D(M_{i+\frac{s}{2}, j+\frac{r}{2}}^{(p)})}{\Delta x^2} \right) - F(C_{i,j}^{(p)}) + \frac{1}{\Delta t} \right). \quad (3.11)$$

<sup>196</sup> This simplifies to,

$$\text{197 } F(C_{i,j}^{(p)}) < \frac{1}{\Delta t} \quad (3.12)$$

<sup>198</sup> The symmetry of  $A$  can be trivially shown if one considers the formation of the diagonals. On a

199 single row, each element corresponds to the adjacent grid points of grid  $i, j$ . As the grid ordering  
 200 counts along, the elements that are equidistance from the diagonal are actually reference to the same  
 201 grid point. Therefore we have symmetry.  $\square$

202 It is important to remark that even though there is a condition for which matrix  $A$  is positive definite  
 203 and symmetric, it realistically will never occur. The condition,  $\frac{1}{F(C)} < \Delta t$ , relates the growth of  
 204 the biomass to the size of timestep selected. Specifically, if a large enough time step is choosen,  
 205 then  $A$  is not guaranteed to converge. When this occurs, it means that a time step, larger then the  
 206 characteristic growth rate of the biomass, has been incorrectly choosen. This means that the there  
 207 would be no relavent results since all the growth, and subsequent reactions, would have occured in a  
 208 single timestep.

209 Given that  $A$  is positive definite and symmetric, the conjugate gradiant method can be used to compute  
 210 the solution. As an added property,  $A$  also happens to be diagonally dominate. This results in  $A$  being  
 211 a M-matrix. It also means that it could be solved using Bi-Conjugate Gradient Method. However  
 212 the Conjugate Gradient method has a faster computation time then Bi-Conjugate Gradiant method for  
 213 this problem and is used for this reason (Barrett et al. (1987)).

214 For solving (3.7), the equation can be rearranged into a quadratic form, substituting in  $G(C)$  from  
 215 (2.15)

$$216 \quad (C^{(p+1)})^2 + \left( \kappa - C^k + \frac{h}{2} \gamma M^{(p+1)} + \frac{h \gamma C^k M^k}{2 \kappa + C^k} \right) C^{(p+1)} + \left( -\kappa C^k + \frac{h \gamma \kappa C^k M^k}{2 \kappa + C^k} \right) = 0. \quad (3.13)$$

217 Using the quadratic equation results in,

$$218 \quad C^{(p+1)} = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a} \quad (3.14)$$

219 for which,

$$a = 1$$

$$\begin{aligned} 220 \quad b &= \kappa - C^k + \frac{h}{2} \gamma M^{(p+1)} + \frac{h}{2} \frac{\gamma C^k M^k}{\kappa + C^k} \\ c &= -\kappa C^k + \frac{h}{2} \frac{\gamma \kappa C^k M^k}{\kappa + C^k} \end{aligned} \quad (3.15)$$

221 Unless  $b^2 - 4ac = 0$ , we have two different solutions to  $C^{(p+1)}$ . The problem with that is that if  
 222 both solutions are positive we have two valid values to be used. If the quantitative behaviour of the  
 223 solutions change with choices of timestep,  $h$ , then the validity of the method goes down. Here, we  
 224 can show that there will always be only one positive solution, regardless of timestep choice.

225 **Proposition 3.2.2.** *The quadratic equation defined as (3.13) will always have one positive solution  
 226 and one negative solution for realistic parameter choices.*

227 *Proof.* Rearranging (3.13) so that all the  $h$  terms are on the right-hand-side, we get

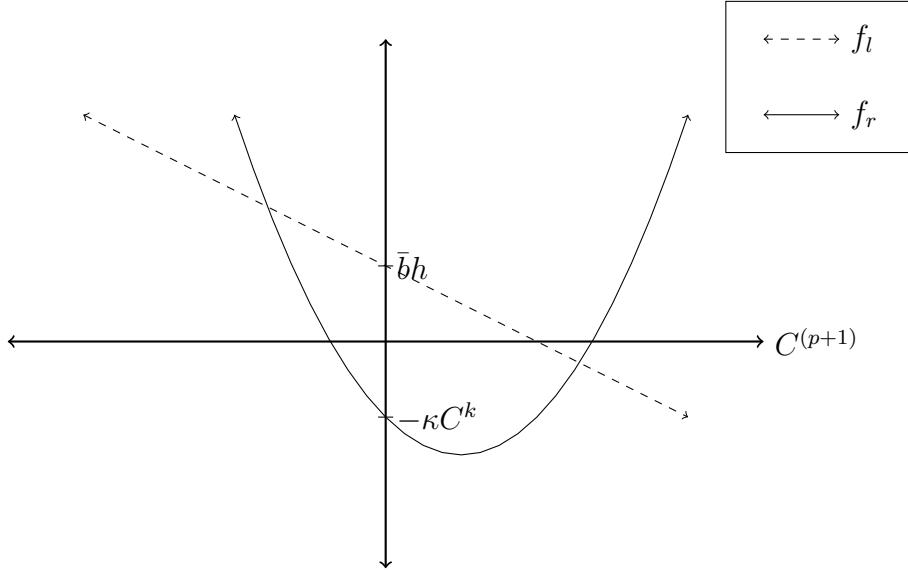
$$228 \quad (C^{(p+1)})^2 + (\kappa - C^k) C^{(p+1)} - \kappa C^k = \left( \frac{\gamma C^k M^k}{2(\kappa + C^k)} - \left( \frac{\gamma M^{(p+1)}}{2} - \frac{\gamma C^k M^k}{2(\kappa + C^k)} \right) C^{(p+1)} \right) h. \quad (3.16)$$

229 To simplify analysis, we let  $\bar{a} = \frac{\gamma M^{(p+1)}}{2} - \frac{\gamma C^k M^k}{2(\kappa + C^k)}$  and  $\bar{b} = \frac{\gamma C^k M^k}{2(\kappa + C^k)}$ .

230 We analyse both the left-hand-side and right-hand-side independently by letting  $f_l = (C^{(p+1)})^2 +$   
 231  $(\kappa - C^k) C^{(p+1)} - \kappa C^k$  and  $f_r = (\bar{b} - \bar{a} C^{(p+1)}) h$ .  $f_l$  is a quadratic equation with positive concavity  
 232 everywhere and  $C^{(p+1)}$ -intercept at  $-\kappa C^k < 0$ .  $f_r$  is a line with a slope opposite to the sign of  $\bar{a}$  and  
 233 has  $C^{(p+1)}$ -intercept at  $\bar{b}h > 0$ . Refer to Figure (3.1) for a visual of  $f_l$  and  $f_r$ .

234 Consider the case  $\bar{a} < 0$ . This means that  $f_r$  has a positive slope. We can show that  $f_r$  must intersect  
 235 with  $f_l$  twice, one at  $C^{(p+1)} < 0$  and another at  $C^{(p+1)}$ .

236 For the positive intersection we can consider a point  $c > 0$ . Because  $f_l = O((C^{(p+1)})^2)$ , using Big-O  
 237 notation, and  $f_r = O(C^{(p+1)})$  we know that there exists a  $c$  for which  $f_l > f_r$  for all  $C^{(p+1)} > c$ .  
 238 Since we trivially have  $f_r(0) > f_l(0)$  we can use the intermediate value theorem to determine that  
 239 there must exist a intersection between  $f_l$  and  $f_r$ . This intersection is the point where  $f_l = f_r$  and



**Figure 3.1:** Graph of  $f_l = (C^{(p+1)})^2 + (\kappa - C^k) C^{(p+1)} - \kappa C^k$  and  $f_r = (\bar{b} - \bar{a} C^{(p+1)}) h$  for  $\bar{a} > 0$  and  $\kappa - C^k < 0$ . Notice that because  $-\kappa C^k < 0$  and  $\bar{b}h > 0$  for all realistic parameter values the two functions will always intersect in the positive  $C^{(p+1)}$  region. Even if  $\bar{a} < 0$  or  $\kappa - C^k > 0$ , the functions are guaranteed to have only one positive intersection since the location of the y-intercepts will not qualitatively change.

240 satisfies our equation resulting in a solution for  $C^{(p+1)} > 0$ . For the negative intersection we know  
 241 that  $f_r(0) > f_l(0)$  and  $f_r$  is monotonically increasing with  $C^{(p+1)}$  and  $f_l$  is monotonically decreasing  
 242 for  $C^{(p+1)} < 0$ . This means that there must exist a point where the two functions intersect and thus  
 243 we have a solution for  $C^{(p+1)} < 0$ .

244 The other case of  $\bar{a} < 0$  follows a similar argument to the previous and also results in one positive  
 245 solution and one negative solution for (3.13).

246 □

247 To determine which branch of (3.14) to use, a physical situation is used. Specifically the case where  
 248 there exist no biomass,  $M = 0$ .

249 The expected outcome is that no substrate is consumed and thus the substrate concentration will  
 250 remain constant as a function of  $t$ . When the equations in (3.15) are evaluated at  $M = 0$ , the result it,

251 
$$a = 1, \quad b = \kappa - C^k, \quad c = -\kappa C^k, \quad (3.17)$$

252 which can be used to evaluate (3.14) as,

$$\begin{aligned}
 C^{(p+1)} &= \frac{-(\kappa - C^k) \pm \sqrt{(\kappa - C^k)^2 - 4(-\kappa C^k)}}{2} \\
 253 \quad &= \frac{1}{2} \left( C^k - \kappa \pm \sqrt{\kappa^2 + 2\kappa C^k + (C^k)^2} \right) \\
 &= \frac{1}{2} (C^k - \kappa \pm (\kappa + C^k)).
 \end{aligned} \tag{3.18}$$

254 Now, if the positive branch is used the above equation evaluates to  $C^{(p+1)} = C^k$ . This means that be-  
 255 tween any two distinct times, the substrate concentration will remain constants, which was expected.  
 256 To further this confirmation, the negative branch results in  $C^{(p+1)} = -\kappa$ , a non-positive substrate  
 257 concentration, which is not physically relavent.

$$258 \quad C^{(p+1)} = \frac{-b + \sqrt{b^2 - 4ac}}{2a} \tag{3.19}$$

259 where  $a$ ,  $b$ , and  $c$  are defined in (3.15).

260 Now that computable solutions for  $M$  and  $C$  at a single time step have been found, an algorithm to  
 261 solve for the next time step can be established. Algorithm 1 shows the organization of solving (3.7 -  
 3.6).

**Data:**  $M^k$ ,  $C^k$  are the values from the previous timestep and  $p = 0$ .

**begin**

```

    Let  $M^{(0)} = M^k$  and  $C^{(0)} = C^k$ ;
    while Convergence is not acheived do
        | Solve  $A^{(p)}M^{(p+1)} = b^{(p)}$ ;
        | Solve  $C^{(p+1)} = \frac{1}{2} (2b \pm \sqrt{b^2 - 4c})$ ;
        | Check convergence;
        | Let  $C^{(p)} = C_{(p+1)}$ ;
        | Let  $M^{(p)} = M_{(p+1)}$ ;
        | Let  $p = p + 1$ ;
    end
    Let  $M^{(k+1)} = M^{(P)}$  and  $C^{(k+1)} = C^{(P)}$ ;
end

```

**Algorithm 1:** Algorithm for the fully-implicit solving of (2.13)

263 Note that Algorithm 1 actually describes both a fully- and semi- implicit method for solving (2.13).  
 264 If  $P = 1$  then only a single iteration of the algorithm is applied, which correlates to a semi-implicit  
 265 method would behave. This would result in a change similar to how the Gauss-Seidal method changes  
 266 the Jacobi method; the values used would no longer be updated in a single timestep when  $P = 1$ .

267 To use the algorithm, the matrix system was converted into a 1D array by use of a bijective mapping  
 268 defined as:

$$\begin{aligned} \pi : \quad \{0, \dots, n\} \times \{0, \dots, m\} &\rightarrow \{0, \dots, nm\} \\ 269 \quad (i, j) &\rightarrow \quad \pi(i, j) \end{aligned} \tag{3.20}$$

270 This mapping allows the system to be easily stored in diagonal format, since (3.10) has five distinct  
 271 diagonals.

### 272 3.3 Computational Setup

273 The implementation of Algorithm 1 was done with Fortran.

274 All the computations were run on a custom built workstation with an Intel Xeon CPU E5-2650 (1.2  
 275 GHz, 20MB cache size) and 32 GB RAM under Red Hat Enterprise Linux Server release 6.5 (Santi-  
 276 ago). Running the computations with OpenMP, took advantage of 4 out of the 16 threads of the Intel  
 277 Xeon CPU, with 2 threads to each core. The selection of 4 threads is because the computation time  
 278 decrease becomes less efficient after more than 4 threads. The GNU Fortran compiler, version 4.4.7,  
 279 was used for all computations; the compiler arguments were

```
280 -O3 -fdefault-real-8 -fopenmp
```

### 281 3.4 Method Validation

282 With a defined method and computational setup a variety of simulations can be run to observe the  
 283 accuracy and behaviour of the method. An examination of a typical simulation will show if the ex-  
 284 pected behaviour is observed, validating the method as functioning. A convergence analysis for the

method can be done to confirm that solutions from different grid sizes approach a single solution as they become more precise. This convergence test will also show the thresholds for an accurate simulation result, to help reduce the computation times. With a well-established method, the comparison between semi- and fully-implicit methods can be done.

### 3.4.1 Basic Simulations

Using Algorithm 1, simple scenarios can be tested as a first verification on the method.

A simple test would be to check if the spatial discretization can preserve specific characteristics of the solutions. One example of this would be seeing if a 1D initial condition could be preserved as time progresses. Having all of the biomass on one boundary of  $\Omega$ , for example across the  $y$ -axis, would qualify as a 1D initial condition. These initial conditions will be defined as:

$$M = \begin{cases} -\left(\frac{h}{d^4}\right)x^4 + h & , \text{if } y \leq d \\ 0 & , \text{otherwise} \end{cases} \quad (3.21)$$

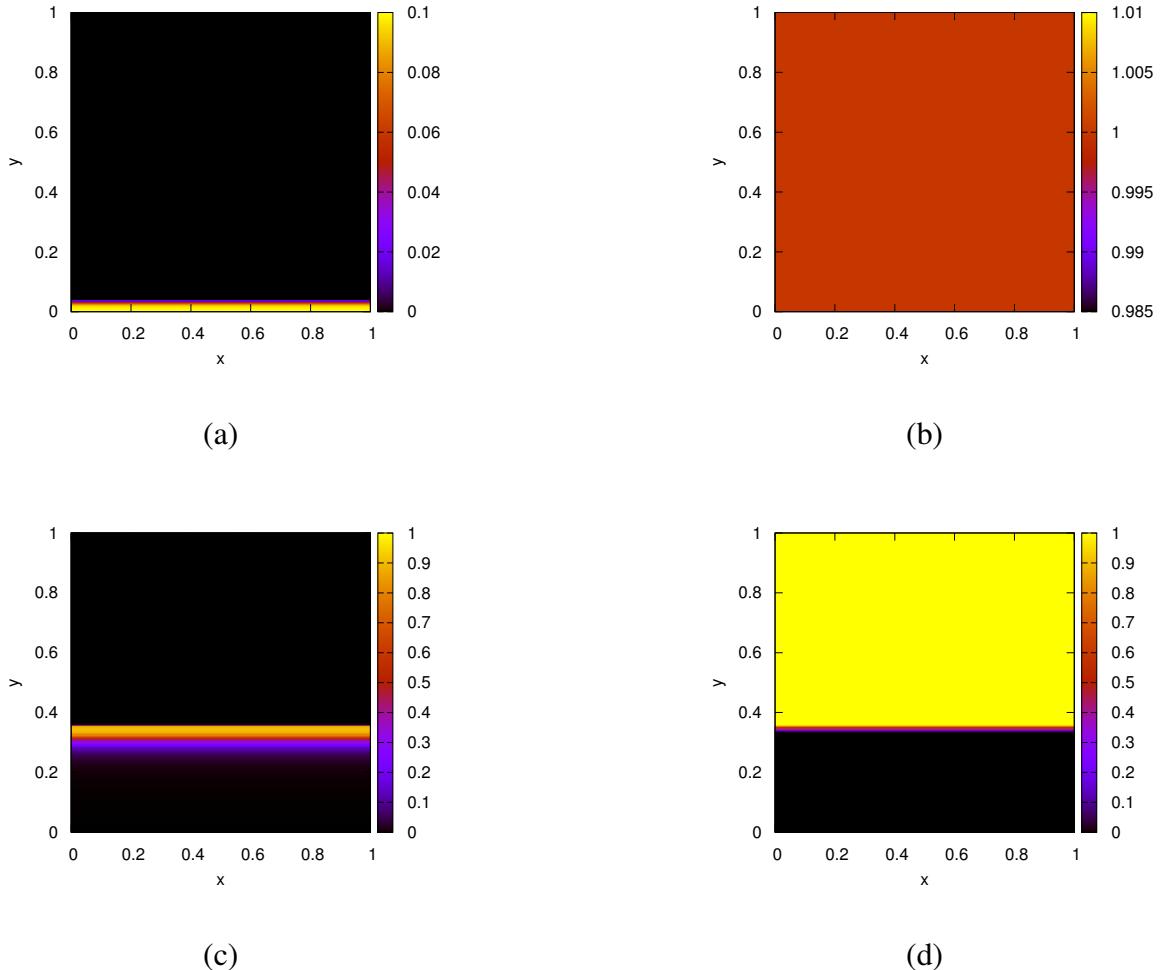
$$C = 1$$

where  $h = 0.1$  and  $d = \frac{5}{128}$ . Here,  $h$  and  $d$  represent the height and depth of the inoculation site.

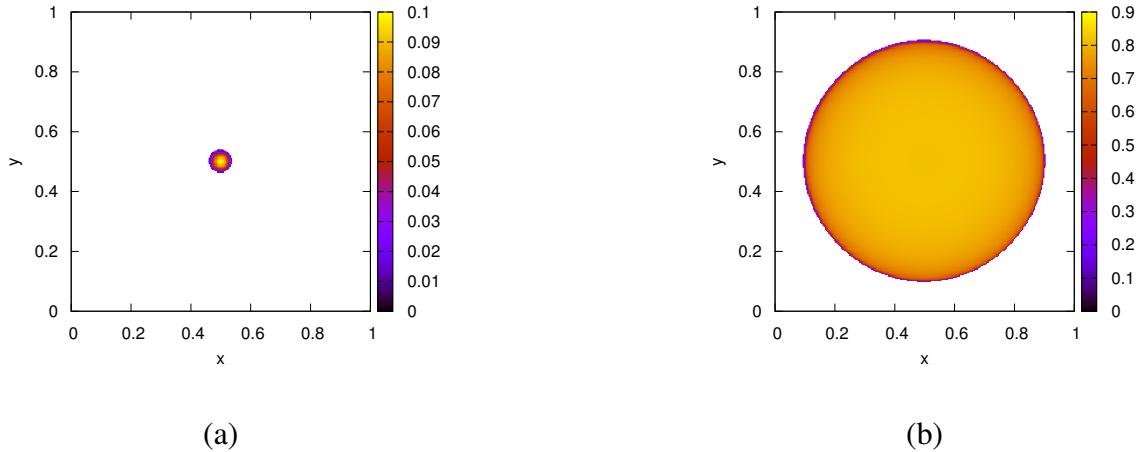
The solution shown in Figure 3.2 shows that the 1D characteristic of the biomass stays at a later time. Another characteristic to observe would be if a spherical initial condition remains spherical. Using initial conditions for the biomass,

$$M = \begin{cases} -\frac{h}{d^2}(x - 0.5)^2 + (y - 0.5)^2 + h & , \text{if } (x - 0.5)^2 + (y - 0.5)^2 < d^2 \\ 0 & , \text{otherwise} \end{cases}, \quad (3.22)$$

a test can be tried to see if the spherical nature of the solution is kept as time progresses. This can be seen in Figure 3.3, where at different times the shape of the biomass,  $M$ , is seen to remain spherical.



**Figure 3.2:** Solutions for (ac)  $M$  and (bd)  $C$  with 1D initial conditions defined in (3.21) at (ab)  $t = 0$  and (cd)  $t = 40$ . Computed with a  $1025 \times 1025$  grid and a timestep of  $\Delta t = 10^{-3}$ .



**Figure 3.3:** Solutions for  $M$  with spherical initial conditions defined by (3.22) at (a)  $t = 0$  and (b)  $t = 40$ . Computed with a  $1025 \times 1025$  grid and a timestep of  $\Delta t = 10^{-3}$ .

303 Both Figure 3.2 and Figure 3.3 increase the confidence that the spatial discretization did not introduce  
304 any loss of characteristics for the solutions.

Given the boundary conditions and spatial discretization, there could be a possible source or sink of biomass when it must diffuse along the boundary of the region. To ensure this is not the case, the total amount of biomass can be used to compare the simulated amount against the theoretical amount.

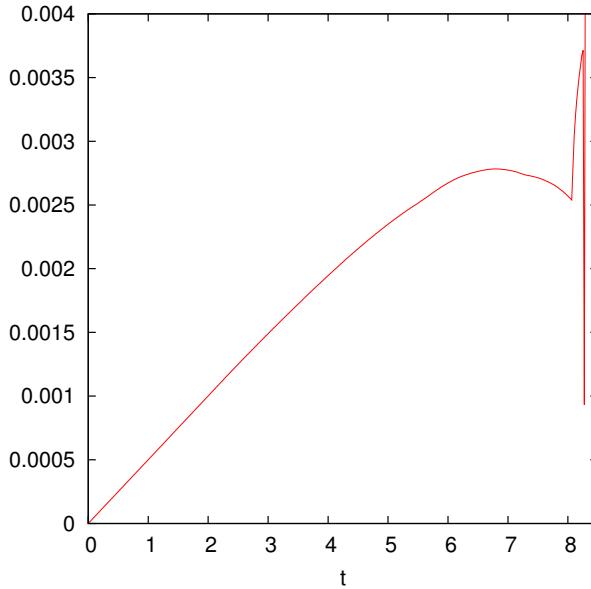
308 However, the total biomass cannot be exactly determined with the given growth rate function. This  
309 means that there will not be anything to measure the validity of the simulation solution against. If we  
310 let the growth rate be some constant called  $a$ , the exact total biomass can be calculated.

<sup>311</sup> The expected total biomass would be of the form  $y_0 e^{at}$ . This can be checked by tracking the total  
<sup>312</sup> biomass, now called  $T_M(t)$ , with the changed growth rate function,  $F(C) = a$ . The calculation of  
<sup>313</sup>  $T_M(t)$  can be done by,

$$T_M(t) = \int_{\Omega} M(t) dx. \quad (3.23)$$

<sup>315</sup> Numerically, this is computed by grid-wise summation,

$$T_M(t^k) \approx T_M^k = \frac{\sum_i^n \sum_j^m M_{i,j}^k}{nm}. \quad (3.24)$$



**Figure 3.4:** Plot of the relative error,  $\frac{|f_1 - f_2|}{|f_2|}$ , between the computed total biomass,  $f_1 = T_M(t)$ , and the theoretical total biomass,  $f_2 = y_0 e^x$ . The changes after  $t = 8$  are from the biomass having completely filled the region  $\Omega$ . This means that there is no physical space for the biomass to occupy and thus the growth slows down to a stop.

317 The simulation setup used will be analogous to that used for Figure 3.3. The one difference will be  
 318 that the simulation here is ran for a longer time to allow the biomass to diffuse along the boundary,  
 319 showing the boundary effects.

320 From Figure 3.4 we can see that the total biomass only differs between the computed value and the  
 321 theoretical value by a relative error less than 0.003. The cases where the error becomes significant are  
 322 from the region being completely filled with biomass, at which point diffusion is no longer possible.  
 323 The error fluctuates violently here because of this. This suggests that the method does not introduce  
 324 any significant sources or sinks of biomass at the boundary of the region.

325 **3.4.2 Convergence Analysis**

326 To validate the accuracy of the method, convergence analyses on the spatial discretizations will need  
 327 to be made. Then the comparison between the semi- and fully-implicit method established in Algo-  
 328 rithm 1 can investigated. First, a metric must be formed to enable consistent comparisons between  
 329 different simulation solutions. This metric will be referred to as the error.

330 **3.4.2.1 Error Computations**

331 The error is computed by taking the relative normed-difference between two solution in the following  
 332 fashion:

$$333 \quad \epsilon_{sol} = \frac{\|u_1 - u_2\|}{\|u_2\|} \quad (3.25)$$

334 where  $u_1$  represents one simulation solution and  $u_2$  represents the solution that is theoretically more  
 335 accurate. The theoretical accuracy of  $u_2$  derives from the fact that most comparisons will be done  
 336 between solutions where one is trivially expected to be more precise. For our purposes, the solutions  
 337 we compare will typically vary in only  $\Delta x$  or between semi- and fully- implicit. These are understood  
 338 to have the relation that a smaller  $\Delta x$ , and that the fully-implicit method with the highest tolerance is  
 339 to be more accurate. There is an assumption that both  $u_1$  and  $u_2$  have the same number of grid points,  
 340 so that the difference can be taken grid-wise.

341 The results of the error computations, named  $\epsilon_{sol}$ , is a numerical value for the difference between two  
 342 solutions. This depends on the norm used during the computations. Here three norms will be used:

$$343 \quad \ell_1 : \|u\|_1 = \frac{1}{nm} \sum_{\pi(i,j)}^{nm} |u_{i,j}| \quad (3.26)$$

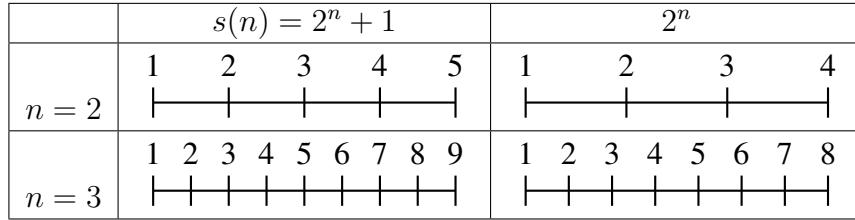
$$344 \quad 345 \quad \ell_2 : \|u\|_2 = \frac{1}{nm} \sqrt{\sum_{\pi(i,j)}^{nm} (u_{i,j})^2} \quad (3.27)$$

$$346 \quad 347 \quad \ell_\infty : \|u\|_\infty = \max_{\substack{i=1, \dots, n \\ j=1, \dots, m}} |u_{i,j}| \quad (3.28)$$

348 These different norms will all be used to create a broader understanding of the error. This creates three  
 349 distinct values for  $\epsilon_{sol}$ , named  $\epsilon_{\ell_1}$ ,  $\epsilon_{\ell_2}$ , and  $\epsilon_{\ell_\infty}$ ; each named for the norm used during the computation.

350 **3.4.2.2 Grid Size Convergence**

351 To observe the validity of the method, a test on the convergence of solutions based on the spatial  
 352 discritization is done. This will involve using the same simulation described in (3.21) due to the



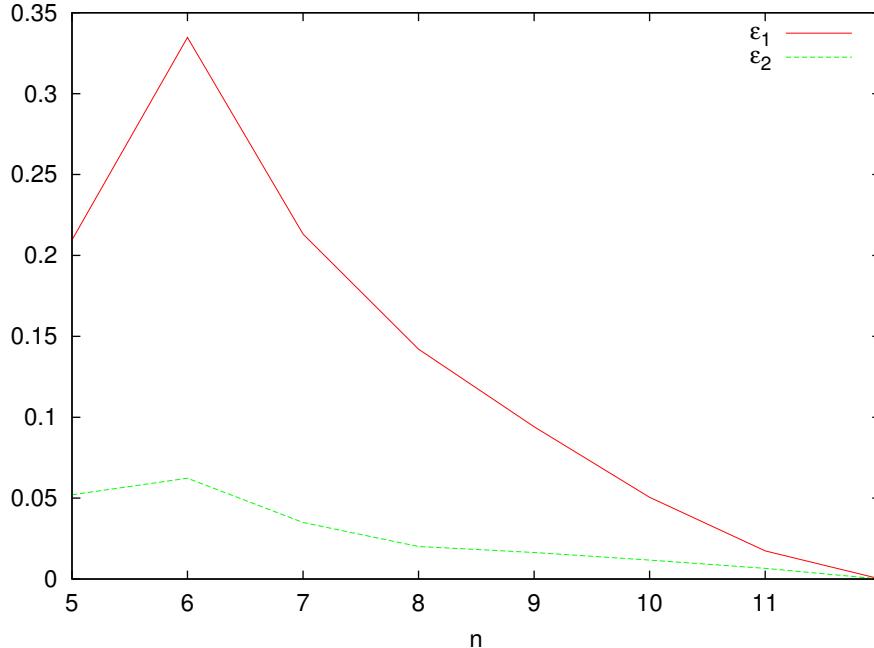
**Figure 3.5:** Visualization in 1D to illustrate the choice of  $s(n) = 2^n + 1$  instead of  $2^n$  for the grid size selection. Here it can be seen that successive grid size selections using  $s(n)$  line up on certain grid points and when using  $2^n$  no grid points are equivalent.

353 simplicity.

354 The convergence will be tracked with only two forms of  $\epsilon_{sol}$ ;  $\epsilon_1$  and  $\epsilon_2$ . This is because the value  
 355 of  $\epsilon_\infty$  doesn't vary with the grid size, since the wave front has a steep interface and tends to lead to  
 356 inconsistent changes in error. Since the use of  $\epsilon_\infty$  is not a suitable method for measuring the error, the  
 357 inconsistency does not suggest an invalidity with the method. Because of the difference in the number  
 358 of grid points between different solutions,  $u_1$  and  $u_2$ , only the grid points in the coarser refinement  
 359 will be used. This places a limitation on the selection of grid-sizes since there must be some grid  
 360 points locations that are the same for two different chosen grid sizes. For this purpose, we define the  
 361 function  $s(n) = 2^n + 1$  for  $n \in \mathbb{N}$  to be used as the grid size selection function. Now certain grid  
 362 points will match without the use of linear interpolation. Here we do not use the typical grid doubling,  
 363  $2^n$ , because no grid points equal between two different grid size selection, as illustrated in Figure 3.5.

364 Using the same simulation setup as was done in Figure (3.2), solutions resulting from different grid  
 365 sizes based on  $s(n)$  are computed for  $n = 5, 6, \dots, 12$ . In this case, when calculating  $\epsilon_{sol} = \frac{|u_1 - u_2|}{|u_2|}$ ,  
 366 we let  $u_1$  be the grid size under investigation and  $u_2$  be the solutions of the most refined gridsize,  
 367  $n = 12$ . This would show the converging solutions for smaller grid sizes because the change to the  
 368 finest grid size will be monotonically decreasing.

369 The results from Figure 3.6 show that the solutions converge as the grid size become refined with the  
 370 grid size. The nonmonotonicity for  $n = 5$  is because the grid has become too large for predictable cal-  
 371 culations, which is an acceptable error since such a coarse grid would never be used for computational  
 372 simulations.



**Figure 3.6:** Plot showing the convergence of solutions based on changes in  $\Delta x$ . The computations are of  $\epsilon_{\ell_1}$  and  $\epsilon_{\ell_2}$  with grid-size following  $s(n) = 2^n + 1$ .

### 373 3.5 Comparison of Semi-implicit and Fully-implicit Method

374 Here the main comparison that analyses the effects of using Algorithm 1 with different tolerances.

375 Recall that the main observation is for  $tol. = 1$ , which correlates to the semi-implicit method since it

376 will allow only a single iteration of the algorithm.

377 The simulation used is the same as described in (3.2). The comparison will be on multiple metrics:

378 the average number of iterations of Algorithm 1, the value of  $\epsilon_1$  and  $\epsilon_2$ , the computation time of the

379 simulation, and the location of the wave peak.

380 The average number of iterations are tracked so that an idea of the extra work can be formed. As the

381 tolerance decreases the amount of iterations the algorithm must perform will increase, the degree of

382 increase will help relate the amount of work.

383 The value of  $\epsilon_1$  and  $\epsilon_2$  act as a measure of accuracy. Here, these values correspond to the difference

384 between a pair of solutions,  $u_1$  and  $u_2$ . The choice of  $u_2$  here is the semi-implicit methods solution.

385 This results in a relative difference from the semi-implicit method and would show the change in solu-

386 tion as the tolerance decreases. Each row of Table 3.1 refers to the  $u_1$  values used in the comparison.

387 Each difference was taken at the last timestep.

388 Along with accuracy, the simulation time is tracked. This is because it represents another metric  
 389 for which the viability of the fully-implicit method can be verified. Theoretically there should be a  
 390 decrease in the error with the fully-implicit method as the value for  $tol$  decreases. Therefore, this  
 391 needs to be weighted against the cost of computational intensity and the increase of the simulation  
 392 time.

393 The location of the wave peak is a tracked quality of the solution that reveals how consistent the  
 394 results are. The wave peak is described here as the maximum value of the solution at the final timestep  
 395 calculated. The ultimate goal is that the simulation solutions be converging towards the exact solution.  
 396 To see this here the  $x$ -coordinate of the wave peak is tracked as well as the height of the wave peak.

The results of the method comparison can be seen in Table 3.1. There are a number of observations

Tol.	Avg. Iter.	$\epsilon_1$	$\epsilon_2$	Time	Wave Height
1.0e-0	1.0000	0.000000000000	0.000000000000	12.1830	0.96366123
1.0e-1	1.0000	0.000000000000	0.000000000000	12.2080	0.96366123
1.0e-2	1.0000	0.000000000428	0.000000000167	12.3379	0.96366123
1.0e-3	1.0000	0.000000000428	0.000000000167	12.2310	0.96366123
1.0e-4	1.0000	0.000000001098	0.000000000329	12.3200	0.96366123
1.0e-5	1.9650	0.002573969658	0.001066499658	18.9869	0.96391491
1.0e-6	2.0000	0.002574057907	0.001066505860	19.0910	0.96391479
1.0e-7	2.0018	0.002573959764	0.001066498736	19.0940	0.96391492
1.0e-8	2.5856	0.002577916965	0.001066781565	20.5169	0.96390966
1.0e-9	2.9012	0.002577461054	0.001066759099	21.3080	0.96390979
1.0e-10	3.2278	0.002581334868	0.001067029069	22.2280	0.96390490
1.0e-11	16.0990	0.002632955188	0.001070709373	57.5589	0.96383105
1.0e-12	36.3184	0.002647234923	0.001071748013	113.9940	0.96380854
1.0e-13	57.6812	0.002648733848	0.001071857564	173.8489	0.96380614

Table 3.1: Results from running simulations with different Tol.

397

398 that can be made from these results.

399 • A positive relationship between computation time and average number of iterations exists.

400 • There is no significant difference between solutions unless the tolerance is set high enough to

401 force multipl iterations. So the semi-implicit method results in a tolerance between  $10^{-4}$  and  
402  $10^{-5}$  since any tolerance forced below that does not require addition iterations.

- 403
- The differences in Wave Height are a result of addition iterations and monotonically approach  
404 a specific value as the tolerance becomes smaller.

405

  - The greatest gain in accuracy while weighing the increased computation time is from a tolerance  
406 around  $10^{-5}$  at which point only one extra iteration is completed.

407 **Chapter 4**

408 **Simulation Results**

409 **References**

- 410 R. Barrett, M. Berry, Chan T.F., J. Demmel, J. Donato, J. Dongarra, V. Eijkhout, R. Pozo, C. Romine,  
411 and H. van der Vorst. *Templates for the Solution of Linear Systems: Building Blocks for Iterative*  
412 *Methods*. Society for Industrial and Applied Mathematics, 1st edition, 1987.
- 413 J.C. Butcher. *Numerical Methods for Ordinary Differential Equations*. Wiley, 2nd edition, June 2008.
- 414 A Dumitrache. *Understanding Biofilms of Anaerobic, Thermophilic and Cellulolytic Bacteria: A*  
415 *Study towards the Advancement of Consolidated Bioprocessing Strategies*. PhD thesis, University  
416 of Toronto, 2014.
- 417 HJ Eberl and L Demaret. A finite difference scheme for a doubly degenerate diffusionreaction equa-  
418 tion arising in microbial ecology. *Electron. J. Differential Equations*, page 7795, 2007.
- 419 HJ Eberl, DF Parker, and MCM van Loosdrecht. A new deterministic spatio-temporal continuum  
420 model for biofilm development. *Journal of Theoretical Medicine*, pages 161–175, 2001.
- 421 S. Geršgorin. über die abgrenzung der eigenwerte einer matrix. *Bulletin de l'Académie des Sciences*  
422 *de l'URSS. Classe des sciences mathématiques et na*, pages 749–754, 1931.
- 423 DR Noguera, G Pizarro, DA Stahl, and BE Rittmann. Simulation of multispecies biofilm development  
424 in three dimensions. *Water Science and Technology*, 39:123–130, 1999.
- 425 BE Rittmann and PL McCarty. Model of steady-state biofilm kinetics. *Biotechnology and Bioengi-  
426 neering*, 22:2343–2357, 1980.
- 427 Y. Saad. *Iterative Methods for Sparse Linear Systems*. Society for Industrial and Applied Mathematic,  
428 2nd edition, 2003.

- 429 S Sirca and M Horvat. *Computational Methods for Physicistsl: Compendium for Students*. Springer,  
430 2012.
- 431 Z-W Wang, S-H Lee, JG Elkins, and JL Morrell-Falvey. Spatial and temporal dynamics of cellulose  
432 degradation and biofilm formation by caldicellulosiruptor obsidiansis and clostridium thermocel-  
433 lum. *AMB Express*, 1:30, 2011.

