# Implementation of 3 stage Lobatto IIIC into the Assimulo package

## Edmund Aristid Lehsten

Bachelor's thesis
2021:K45

## Lund University

Faculty of Science
Centre for Mathematical Sciences
Numerical Analysis

# Implementation of 3 stage Lobatto IIIC into the Assimulo package

Edmund Aristid Lehsten

2021

# Abstract

In recent years, the popularity of discontinuous Galerkin methods has increased. As shown in [19], a result exists that states that the Discontinuous Galerkin space approximations (DG) are equivalent to the Lobatto IIIC Runge-Kutta method. This thesis therefore outlines the adaptation of Hairer's implementation of the Radau IIA Runge-Kutta method to the Lobatto IIIC method, extended with an adaptation of Pinto et al.'s two step error estimation found in [17]. As an alternative to the classical Three staged Radau IIA Runge Kutta method.

# Popular Scientific Description

Ordinary differential equations (ODEs) are the backbone of physics, chemistry and biology. However these cannot be solved over a continuous region in time and space, rather they need to be solved on a discrete grid, while iterating over different points in time.

Runge Kutta methods are one of the methods to solve these problems. In this thesis we discuss the implementation of one such method, the Lobatto IIIC, by adapting the implementation of a similar method, the Radau IIA, done by Hairer [14]. Further we expand upon the implementation by including a newer method by pinto et al. for the local error estimation at the different steps [17].

# Acknowledgements

# Introduction

In numerical mathematics a common field of interest is solving ordinary differential equations. These equations are commonly used when describing how values that depend on each other change over time. Many different methods for solving these problems exist, one group being the family of Runge-Kutta methods. However, it is not possible to solve these problems, for example the temperature in a room, at every single point within that region. Therefore the regions get discritized in space, which means that one picks a discrete number of points within this region and reformulates the problem to describe the relation between points. Another type of discretization is temporal discretization, where instead of picking points in space one picks points in time.

In recent years, discontinuous Galerkin methods have become more popular than previousely [19]. These methods allow one to discritize in both space and time. This is benificial since it allows us to prove different properties mathematically, such as, for example entropy stability. However these methods are not as easy to use as the methods that we have for solving ODEs. Luckily a result exists proving that the Lobatto IIIC Runge Kutta method is equivalent to one of hte Discontinuous Galerkin methods, specifically the spatial Discontinuous Galerkin approximation [19]. Hence when we use the Lobatto IIIC method we have the same properties as we have for the Discontinuous Galerkin space approximation. There has not been to much focus on the Lobatto IIIC method since it has lower order than the Radau method. Furthermore the difference between the methods is solely in their coefficient. This work therefore aims to adapt Heirer's implementation for the Radau IIA Runge Kutta method into an implementation of the Lobatto IIIC method [14]. Furthermore we will incorporate newer results from pinto et al. that improve upon the error estimation disccused by Hairer [17].

We begin by giving a formal definition for the Runge-Kutta method in Section 2. Section 2.1 focuses on the initial conversion of the method into a linear algebra setting. In Section 2.2 we focus on the Newton iteration in the method, along with its starting and stopping conditions. Section 2.3 discusses further simplifications that can be done to the method to improve simulation time. In Sections 2.4 and 2.5 we discuss the approximation of the error and how to apply this to dynamically adjust the step size of the method. We then mention a few differences between the coefficients and their corresponding Runge-Kutta methods while also introducing the Lobatto IIIC coefficient in section 2.6
Section 3 focuses on the package we implemented our method in, while Section 4 discusses the numerical results.

# Contents

# Chapter 1

# Prerequisites

Throughout this work we assume general knowledge of numerical linear algebra and numerical analysis. However, in the rest of this section we shall reiterate some of the concepts that will be used heavily throughout the rest of this work.

## 1.1  Numerical Linear Algebra

We shall begin this section by giving the definition of the Kronecker product which we utilize later throughout this work and then follow up with the LU-decomposition in conjunction with the Newton iteration.

**Definition 1.1.1.** *The **Kronecker product** ($\otimes$) between two matricise $\boldsymbol{A}$ and $\boldsymbol{B}$ is defined as*

$$\boldsymbol{A} \otimes \boldsymbol{B} = \begin{bmatrix} a_{1,1}\boldsymbol{B} & \cdots & a_{1,n}\boldsymbol{B} \\ \vdots & \ddots & \vdots \\ a_{m,1}\boldsymbol{B} & \cdots & a_{m,n}\boldsymbol{B} \end{bmatrix}$$

Given a function $f$, $f : \mathbb{R}^n \to \mathbb{R}^n$, the Newton iteration is a method which attempts to find zeros of the function $f$. In other words, we obtain an $x$ such that $f(x) = 0$. This method is defined by the following iterative formula.

$$x^{(k+1)} = x^{(k)} - \frac{f(x^{(k)})}{f'(x^{(k)})}$$

**Remark 1.1.2.** *The Newton iteration can be generalized to multiple dimensions by using the Jacobian $\mathscr{J}$ of $f$ instead of $f'$ and solving the following instead iteratively*

$$\mathscr{J} \Delta x^{(k)} = f(x^{(k)}) \tag{1.1}$$

$$x^{(k+1)} = x^{(k)} - \Delta x^{(k)} \tag{1.2}$$

Note the equation (1.1) is a linear system of equations.

If the starting point, $x^{(0)}$, of the Newton iteration is close enough to the wanted solution, we can use something called the **simplified Newton iteration**. The difference between these two

methods is that the Jacobian (or derivative) is only evaluated once and then reused for the rest of the steps. The benefit with this is that computing the Jacobian can be costly depending on the function, the simplified method removes this cost as it only computes it once. Despite that the simplified Newton method comes with the trade of that it is only locally linearly convergent instead of being quadratically convergent like the classical method. This means that the domain in which a starting point will converge to a zero can be smaller depending on the function. Further more the simplified Newton method often requires slightly more iterations to converge. However since the computation of the Jacobian can be extremely costly overall the simplified Newton method is still an improvement [7]. To further improve upon this method we introduce the LU-decomposition

**Definition 1.1.3.** *An **LU-decomposition** of a square matrix A is the decomposition $A = LU$, where L is a lower triangular matrix and U is an upper triangular matrix.*

The benefit of the LU-decomposition becomes apparent when repeatedly solving linear systems of equations. Solving a lower/upper triangular system is done by either forward or backwards substitution and hence is faster to compute than a full system. When solving $Ax = b$ we instead solve for $Ly = b$ and then for $Ux = y$. Both of these are straight forward to compute and hence solving a system given the LU-decomposition of $A$ is more efficient than solving the original system. Therefore, since we want to repeatedly solve $\mathcal{J}\Delta x^{(k)} = f(x^{(k)})$, for different $x^{(k)}$, we can compute the LU-decomposition of $\mathcal{J}$ once and use it for the repeated Newton iterations to make the computations more efficient.

## 1.2 Ordinary Differential Equations and Numerical Mathematics

*This section is based no chapter 12 in [18].*
Ordinary differential equations (ODE)s are common in numerical mathematics. Many different methods exist for finding approximate solutions to these problems. The aim of this work is to implement such a method. However, before we can begin discussing the method we will focus on later, some apriori knowledge is required. We begin with the definition of an ODE and initial value problem.

**Definition 1.2.1.** *An **ordinary differential equation** (ODE) is an equation of the form*

$$y'(x) = f(x, y(x)) \tag{1.3}$$

*where $y(x)$ is a function and $x$ is a variable, commonly $x$ is considered to be time.*

As with integrating where the solution contains an arbitrary constant, ODEs can have an infinite amount of solutions known as **solution curves**.

To uniquely define just one solution curve, an initial value $(x_0, y_0)$, with $y(x_0) = y_0$, is given. This value is a form of starting point from which we begin drawing our curve.

**Definition 1.2.2.** *An **initial value problem** is an ODE together with an initial value.*

Commonly a numerical method can only find an approximate solution $\tilde{y}(x)$ of $y(x)$ for a given $x$. To arrive at a certain point $x$, most methods take multiple steps calculating $y_n \approx y(x_n)$ for each step, i.e., $x_0 \rightarrow x_1 \rightarrow x_2 \rightarrow \cdots \rightarrow x_n = x$. The step size of a step $x_i \rightarrow x_{i+1}$ is $h_i := x_{i+1} - x_i$.

To measure the accuracy of these methods we can observe the **global error** of the method. That is we can chooses a point $x_n$ and compute $\tilde{y}(x_n)$ and then compare this result with the solution $y(x_n)$. However, in most cases a true solution in unknown. We can therefore not evaluate the accuracy of the method. To ensure that the method is still sufficiently accurate we instead study a **local error**[*] estimate which compares two methods of different order at the same point, as found in [8]. We denote this error at $x_n$ as $T_n$: the greater $T_n$ is the greater the local error is. The order of a method is defined as follows

**Definition 1.2.3.** *A numerical method is said to have **order (of accuracy)** p, if p is the largest positive integer such that, for any sufficiently smooth solution curve $(x, y(x))$ in the solution domain of the initial value problem, there exist constants $K$ and $h_0$ such that*

$$|T_n| \leq Kh^p \quad for \ 0 < h \leq h_0$$

*for any pair of points $(x_n, y(x_n))$, $(x_{n+1}, y(x_{n+1}))$ on the solution curve*

**Definition 1.2.4.** *A numerical method is **consistent** with the differential equation* (1.3) *if the local error, is such that for any $\epsilon > 0$ there exists a positive $h(\epsilon)$ for which $|T_n| < \epsilon$ for $0 < h < h(\epsilon)$ and any pair of points $(x_n, y_n), (x_{n+1}, y_{n+1})$ on the solution curve.*

---

[*]Also known as the Local Truncated Error

# Chapter 2

# The Runge-Kutta Method

The idea of the Runge-Kutta method is to extend upon the standard Euler method, as defined in [18], for numerical integration. Instead of evaluating the function $f$ once per step, one evaluates it multiple times. This idea of using a multiplicity of evaluations was originally proposed by Runge in 1895, and was developed further by Heun in 1900 and Kutta in 1901. The latter characterized the set of Runge-Kutta methods of order 4 and proposed the first methods of order 5, as Butcher points out in [8].

Given an initial value problem as defined in definition 1.2.2, we can use the Runge-Kutta method to take a step from $x_0 \to x_1 = x_0 + h$ and find $y_1 = f(x_1)$.

**Definition 2.0.1.** *An s-stage Runge-Kutta method is defined as:*

$$g_i = y_0 + h \sum_{j=1}^{s} a_{i,j} f(x_0 + c_j h, g_j) \quad , i = 1, \ldots, s \tag{2.1a}$$

$$y_1 = y_0 + h \sum_{j=1}^{s} b_j f(x_0 + c_j h, g_j) \tag{2.1b}$$

*with $a_{i,j}, b_i, c_i$ being constant coefficients that are method specific (Radau IIA, Lobbato IIIC, etc.), and h being the step size .*

Note that $h$ needs to be sufficiently small for the Runge-Kutta method to be stable.

The Runge-Kutta methods can be roughly grouped into two variants, the explicit and the implicit methods, the explicit methods having the further restriction on the coefficient that $a_{i,j} = 0 \quad \forall j \geq i$, this alters the sum in (2.1a) and the method becomes as follows:

**Definition 2.0.2.** *An s-stage explicit Runge-Kutta method is defined as:*

$$g_i = y_0 + h \sum_{j=1}^{i-1} a_{i,j} f(x_0 + c_j h, g_j), \quad i = 1, \ldots, s \tag{2.2a}$$

$$y_1 = y_0 + h \sum_{j=1}^{s} b_j f(x_0 + c_j h, g_j) \tag{2.2b}$$

*with $a_{i,j}, b_i, c_i$ being constant coefficients that are method specific.*

We are interested in the Lobatto IIIC method which is an implicit method. We therefore refer to definition 2.0.1 when we talk about the a Runge-Kutta method throughout the rest of this work.

The selection of these constants results from further conditions imposed on the method ensuring it is consistent and converges. However, the creation of these coefficients is beyond the scope of this text. For further reference refer to Butcher [8, section 33 & 340]. Additionally, Butcher introduced the following structure for grouping the coefficients [8]

$$
\begin{array}{c|cccc}
c_1 & a_{1,1} & a_{1,2} & \cdots & a_{1,s} \\
c_2 & a_{2,1} & a_{2,2} & \cdots & a_{2,s} \\
\vdots & \vdots & \vdots & & \vdots \\
c_s & a_{s,1} & a_{s,2} & \cdots & a_{s,s} \\
\hline
& b_1 & b_2 & \cdots & b_s
\end{array}
$$

For further convenience we shall group the individual coefficients into a matrix and vectors following the same structure:

$$
A = \begin{bmatrix}
a_{1,1} & a_{1,2} & \cdots & a_{1,s} \\
a_{2,1} & a_{2,2} & \cdots & a_{2,s} \\
\vdots & \vdots & & \vdots \\
a_{s,1} & a_{s,2} & \cdots & a_{s,s}
\end{bmatrix}
$$
$$
b = (b_1, \cdots, b_s)
$$
$$
c = (c_1, \cdots, c_s)
$$

## 2.1   Initial Implementation

*The following sections are based on chapter IV.8 from Hairer [14] unless stated otherwise.*
To start the implementation of the implicit Runge-Kutta method, we begin by introducing $z_i$ as

$$
z_i = g_i - y_0 \tag{2.3}
$$

This change is done to reduce round-off errors in floating point calculations[*]. $g_i$ are the stage values, i.e., the values that we would have at times $t + c_i h$ respectively. By looking at the difference we get that $z_i$ are the increments of the stage values. This is desired since the increments should be small when taking small steps. Observe that the terms in (2.3) are all vectors of dimension $n$. By calculating the increment instead of the actual value we avoid round-off errors.

Incorporating these terms into equation (2.1a) we get

$$
z_i = h \sum_{j=1}^{s} a_{i,j} f(x_0 + c_j h, y_0 + z_i) \ \ i = 1, \dots, s \tag{2.4}
$$

---

[*]When considering all numbers that can be represented by the float point data structure, these numbers are more dense around 0 and less the further away they are.

Note equation (2.4) can be written in vector notation as follows:

$$\mathbf{Z} = \begin{pmatrix} z_0 \\ z_1 \\ \vdots \\ z_s \end{pmatrix} = A \begin{pmatrix} hf(x_0 + c_1 h, y_0 + z_1) \\ hf(x_0 + c_2 h, y_0 + z_2) \\ \vdots \\ hf(x_0 + c_s h, y_0 + z_s) \end{pmatrix} \tag{2.5}$$

Observe that each term $z_i$ is a vector of size $n$.

Then, assuming $A^{-1}$ exists, we can introduce $(d_1, \ldots, d_s) = (b_1, \ldots, b_s)A^{-1}$, and equation (2.1b) becomes:

$$y_1 = y_0 + \sum_{j=1}^{s} d_j z_j \tag{2.6}$$

Note that in equation (2.6) the $d_j$ are scalars. Furthermore the $d_j$'s are only defined by the method constants and can thus be pre-calculated. This removes the need to determine the inverse of $A$ while using the method.

## 2.2 Newton's Method

We now use Newton's method to evaluate $\vec{Z}$ in equation (2.5). For this we rearrange the equation to:

$$\vec{Z} - A \begin{pmatrix} hf(x_0 + c_1 h, y_0 + z_1) \\ hf(x_0 + c_2 h, y_0 + z_2) \\ \vdots \\ hf(x_0 + c_s h, y_0 + z_s) \end{pmatrix} = 0 \tag{2.7}$$

The Jacobian of this system is:

$$I - h \begin{pmatrix} a_{1,1} \frac{\delta f}{\delta y}(x_0 + c_1 h, y_0 + z_1) & \cdots & a_{1,s} \frac{\delta f}{\delta y}(x_0 + c_s h, y_0 + z_s) \\ & & \\ a_{s,1} \frac{\delta f}{\delta y}(x_0 + c_1 h, y_0 + z_1) & \cdots & a_{s,s} \frac{\delta f}{\delta y}(x_0 + c_s h, y_0 + z_s) \end{pmatrix} \tag{2.8}$$

Note that in equation (2.8) each term $\frac{\delta f}{\delta y}(\cdot, \cdot)$ is a matrix of dimension $n \times n$. Moreover these terms can all be approximated as follows:

$$\frac{\delta f}{\delta y}(x_0 + c_i h, y_0 + z_i) \approx \frac{\delta f}{\delta y}(x_0, y_0) = \mathscr{J}$$

This allows us write an approximation to (2.8) as follows:

$$I - hA \otimes \mathscr{J} \tag{2.9}$$

where $\otimes$ is the Kronecker product as in definition 1.1.1. Note that from this point on $I$ always denotes the identity matrix of appropriate dimension.

Thus, the Newton iteration becomes

$$(I - hA \otimes \mathscr{J})\Delta\vec{Z}^k = -\vec{Z}^k + h(A \otimes I)F(\vec{Z}^k) \tag{2.10a}$$

$$\vec{Z}^{k+1} = \vec{Z}^k + \Delta\vec{Z} \tag{2.10b}$$

with $F(\vec{Z}^k) = \big(f(z_0), \ldots, f(z_s)\big)$

### 2.2.1 Newton Starting Value

There are different alternatives for the starting value of the Newton iteration. Since we are taking a small step and are calculating the increments, we know that they should be small. Hence setting $Z^0 = \mathbf{0}$ (that is the zero vector of suitable size) is a valid choice, and is what we decided to use. However, other options exist, for instance one could take a small initial step using the explicit Euler method and then use that as the starting value. These other methods have the benefit of arriving at the stopping condition (see next section) in fewer iterations with the downside of being slightly more computationally intensive in the setup. However, the reduction in iterations can often compensate for their initial setup expense and, overall, these methods could yield faster solution compared to the basic approach. Our choice was made to simplify the implementation and since the final problem we will test this method on, see section 4, is linear and a single iteration will be sufficient.

### 2.2.2 Newton Stopping Condition

For the stopping condition to the Newton iteration we again look to Hairer [14, see IV.8]. We begin by estimating the convergence rate

$$\Theta_k = \frac{||\Delta Z^k||}{||\Delta Z^{k-1}||} \tag{2.11}$$

where $\Delta Z^k$ is the result of this iteration and $\Delta Z^{k-1}$ the result of the previous iteration. Then, we define the stopping condition to stop iterating if

$$\eta_k ||\Delta Z^k|| \leq \kappa \cdot Tol \quad \text{with} \quad \eta_k = \frac{\Theta_k}{1 - \Theta_k} \tag{2.12}$$

holds[†].

Once stopped we accept $Z^k$ as the approximation $Z^*$ to $Z$. Here $Tol$ denotes the desired tolerance of the solution. As mentioned by Hairer, $\kappa$ is arbitrary and must be chosen by the user, Hairer suggest a value around $10^{-1}$ or $10^{-2}$. To reduce the necessary Jacobian computation, we shall only recompute the Jacobian if the Newton iterations start to diverge, which we define by $\Theta_k \geq 1$.

## 2.3 Simplification for Implementation

The Runge-Kutta method discussed above has a few drawbacks, mainly in the computational intensity. Currently for each Newton iteration we must compute:

- $s$ evaluations of $f$

- solve an $n \cdot s \times n \cdot s$-dimensional linear system

---

[†]Note that this stopping condition for the Newton iteration requires a minimum of two steps. Hairer makes a suggestion on how this could be improved for linear (or near linear) problems where one iteration would be sufficient. We have implemented it in code but left it out of this work as it is of no further interest.

If the diagonalization of $A^{-1}$ exists i.e.

$$T^{-1}A^{-1}T = \Lambda \tag{2.13}$$

then, by pre-multiplying (2.10a) with $(hA^{-1}) \otimes I$ and introducing the term $W^k = (T^{-1} \otimes I)Z^k$, equations (2.10a) and (2.10b) become:

$$(H^{-1}\Lambda \otimes I - I \otimes J)\Delta W^k = -h^{-1}(\Lambda \otimes I)W^k + (T^{-1} \otimes I)F\big((T \otimes I)W^k\big) \tag{2.14a}$$

$$W^{k+1} = W^k + \Delta W^K \tag{2.14b}$$

From this point on we will focus mainly on the Runge-Kutta methods with $s = 3$[‡].
For the ideal case that $A^{-1}$ has three real eigenvalues $\lambda_1, \lambda_2, \lambda_3$ then

$$\Lambda = \begin{bmatrix} \lambda_1 & 0 & 0 \\ 0 & \lambda_2 & 0 \\ 0 & 0 & \lambda_3 \end{bmatrix} \tag{2.15}$$

This would imply that the matrix in front of $\Delta W^k$ in equation (2.14a) would become:

$$\begin{bmatrix} \lambda_1 I - \mathscr{J} & 0 & 0 \\ 0 & \lambda_2 I - \mathscr{J} & 0 \\ 0 & 0 & \lambda_3 I - \mathscr{J} \end{bmatrix} \tag{2.16}$$

Note here that each element is a block matrix of size $n \times n$. Solving this system can be simplified by splitting it into solving three $n \times n$ systems instead of solving one $3n \times 3n$ system. Unfortunately, most for most Runge-Kutta methods their respective $A^{-1}$ matrices do not have three real eigenvalues. Rather they have one real, $\gamma$, and a complex pair, $\alpha \pm \beta i$. In this case, the matrix infront of $\Delta W^k$ in equation (2.14a), (the equivalent to (2.16)) would look as follows:

$$\begin{bmatrix} \gamma I - \mathscr{J} & 0 & 0 \\ 0 & \alpha I - \mathscr{J} & -\beta I \\ 0 & \beta I & \alpha I - \mathscr{J} \end{bmatrix} \tag{2.17}$$

Here again one can split it into smaller sections to help the computer while solving this. Therefore instead of one $3n \times 3n$ system it would come down to one $n \times n$ and one $2n \times 2n$ system. For the $2n \times 2n$ system further tricks exist to aid the computer in solving this. However, that is beyond the scope of this text[§]. Since the left hand side (LHS) of equation (2.14a) does not change throughout the consecutive iterations within one time step we compute the LU-decomposition of these two matrices (the $n \times n$ and $2n \times 2n$) to significantly decrease the solve time of the systems. If the function and Jacobian are relatively easy to evaluate then this decomposition becomes the bottle neck, that is to say the most costly part, of the method. Hence, if the step size does not change one can reuse the decomposition from the previous step.

---

[‡]This is done since the source text [14] also makes this choice.
[§]for reference refer to Hairer page 122 [14]

## 2.4 Local Truncation Error

To estimate the local error we use a similar method to Hairer in chapter IV.8 [14]. His method consists of computing an approximate solution $\hat{y}_n \approx y_n$ which is of slightly lower order. Then the difference between these two solutions gives an estimate of the error $err = \hat{y}_n - y_n$.

$$\hat{y}_{n+1} = y_n + h \sum_{i=1}^{3} \hat{b}_i z_i \tag{2.18}$$

**Remark 2.4.1.** *These error methods are also known as embedded methods since they reuse the $z_i$ terms from the current step.*

It now remains to find coefficients $\hat{b}_i$ such that the method in (2.18) is of the desired order. As mentioned by Pinto et al. in [17] finding embedded methods of order higher than $s-1$ is in most cases not possible. This means that our lower order approximation will be of order 2. One possible choice would be $(\hat{b}_1, \hat{b}_2, \hat{b}_3) = \left(-\frac{1}{2}, 2, -\frac{1}{2}\right)$ as found in [20]. On the other hand the difference between the order of the two methods is not ideal. We therefore turn to Butcher's idea found in [8], where he suggests to look for a method with $s+1$ stages where the first $s$ stages are identical. With this we would hope to find a method of order 3. Unfortunately, as mentioned, it is unlikely to find such embedded methods for methods of order $\geq 2s-2$. Furthermore, these methods would require more function evaluations which is undesired since they could be costly. Hairer instead suggests using the method of order $s-1$ but adding an explicit step at it and arrive at the following method:

$$\hat{y}_{n+1} = y_n + h \hat{b}_0 f(x_0, y_0) + \sum_{i=1}^{3} \hat{b}_i z_i \tag{2.19}$$

Hairer further suggests the choice of $\hat{b}_0 = \gamma$, with $\gamma$ being the real eigenvalue of $A^{-1}$. Through this explicit step the hope is to improve the order up to $s$. In our case, $s = 3$, the embedded method would then have order 3 while the actual method has order 4. Note, however, that the embedded method still requires a further function evaluation. Consequently, we look at the suggestion by Fabien and Jason in [13]. Their idea was that instead of one step of size $h$, one makes two of size $h/2$. Then an approximation can be computed as

$$\hat{y}_{n+1} = y_n + \sum_{i=1}^{2s} \hat{b}_i z_i \tag{2.20}$$

where $z_1, \ldots, z_s$ are the stage values for the first step and $z_{s+1}, \ldots, z_{2s}$ are the stage values for the second step. This method has one drawback which is that whenever a step is rejected one must reject the last two steps since they must always be of equal step size. This issue was addressed by Pito et al. [17]. Their idea was to adjust the coefficients $b_i$ depending on the ratio between the last two step sizes.

Assume that the last step $y_{n-1} \to y_n$ with $t_n = t_{n-1} + h_{n-1}$ has been successful. Then for the approximate solution of the next step $y_n \to y_{n+1}$ with $t_{n+1} = t_n + h_n$, we have

$$\hat{y}_{n+1} = y_n + \sum_{i=1}^{3} \delta_{n,i} z_{n-1,i} + \sum_{i=1}^{3} \beta_{n,i} z_{n,i} \tag{2.21}$$

here $z_{n-1,i}, z_{n+1,i}$ denote the stage values from $y_{n-1} \to y_n$ and $y_n \to y_{n+1}$ respectively and $\delta_{n,i}, \beta_{n,i}$ both depend on $r_n = \frac{h_n}{h_{n-1}}$. These coefficients can be determined from the following conditions:

$$\beta_n^T e_s = 0, \ \delta_n^T \mathbf{1} = 0 \tag{2.22}$$

$$\delta_n^T A (c-\mathbf{1})^{j-1} + r_n^j \beta_n^T A c^{j-1} = \frac{r^j}{j} \quad j = 1, \dots, \hat{p} \tag{2.23}$$

where $e_s$ is a vector of dimension $s$ with all zeros except for the last term which is 1, $\mathbf{1}$ is an $s$-dimensional vector filled with 1's and $\hat{p}$ is the desired approximation order. Note that for our case, $s = 3$ an $\hat{p} = 3$, equations (2.22) and (2.23) give us and under determined system. We therefore wish to impose a further condition, to that end we must define the stability function $r(z)$ of a numerical method.

*This section on the stability function is based on [15]*
We being by defining the linear stability domain of a method.

**Definition 2.4.2.** *For a numerical method with constant step size $h$, with results $y_n \approx y(x_0 + h \cdot n)$ applied to the linear ODE*

$$y'(t) = \lambda y(t), \ \ t \geq 1, \quad y(0) = 1, \tag{2.24}$$

*the **stability domain** $\mathscr{D}$, of the method, is the set of all numbers $h \cdot \lambda \in \mathbb{C}$ such that the limit $\lim_{n \to \infty} y_n = 0$.*

**Remark 2.4.3.** *Note that for $Re(\lambda) < 0$ we have that for the exact solution of (2.24), the limit $\lim_{t \to \infty} y(t) = 0$*

When applying a Runge-Kutta method to (2.24), we get that

$$\xi_i = y_0 + h \sum_{j=1}^{s} a_{i,j} \xi_j, \quad i = 1, \dots, s, \tag{2.25a}$$

$$y_{n+1} = y_n + h \lambda \sum_{j=1}^{s} b_j \xi_j, \tag{2.25b}$$

with

$$\xi = \begin{bmatrix} \xi_1 \\ \xi_2 \\ \vdots \\ \xi_s \end{bmatrix} = (I - h \lambda A)^{-1} \mathbf{1} y_n, \tag{2.26}$$

where $\mathbf{1}$ is the vector of dimension $s$ with all ones. We denote by $\mathbb{P}_\alpha$ the set of polynomials of maximum order $\alpha$. Similarly, by $\mathbb{P}_{\alpha/\beta}$ we denote the set of rational functions $\hat{p}/\hat{q}$ where $\hat{p} \in \mathbb{P}_\alpha$ and $\hat{p} \in \mathbb{P}_\beta$.

**Lemma 2.4.4.** *For every Runge-Kutta method of order $s$ applied to (2.24) there exists $\rho \in \mathbb{P}_{s/s}$ such that*

$$y_n = \left[ \rho(h \lambda) \right]^n, \quad n = 0, 1, \dots \tag{2.27}$$

*Proof. We shall only give a brief outline of the proof here.*
From (2.25b),(2.26) and (2.27) we get that

$$
\begin{aligned}
y_{n+1} &= y_n + h\lambda \sum_{j=1}^{s} b_j \xi_j \\
&= \left[1 + h\lambda b^{\mathrm{T}}(I - h\lambda A)^{-1}\mathbf{1}\right] y_n \\
\Rightarrow \left[\rho(h\lambda)\right]^{n+1} &= \left[1 + h\lambda b^{\mathrm{T}}(I - h\lambda A)^{-1}\mathbf{1}\right] \cdot \left[\rho(h\lambda)\right]^n \\
\Rightarrow \rho(h\lambda) &= 1 + h\lambda b^{\mathrm{T}}(I - h\lambda A)^{-1}\mathbf{1} \\
\Rightarrow \rho(z) &= 1 + z b^{\mathrm{T}}(I - zA)^{-1}\mathbf{1}, \quad \text{with } z \in \mathbb{C}
\end{aligned}
\tag{2.28}
$$

The rest of this proof now consists of proving that $\rho(z)$ is indeed an element of $\mathbb{P}_{s/s}$, the proof of which we shall not shown here but can be found in [15]. $\qquad\square$

**Lemma 2.4.5.** *Suppose that an application of a numerical method to the linear ODE* (2.24) *produces a solution sequence* $y_n = \left[\rho(h\lambda)\right]^n$ *,* $n = 0, 1, \ldots$*, where $\rho$ is an arbitrary function. Then*

$$
\mathscr{D} = \{z \in \mathbb{C} : |\rho(z)| < 1\}.
\tag{2.29}
$$

*Proof.* The proof of this follows directly from the definition of $\mathscr{D}$ in Definition 2.4.2 $\qquad\square$

As Pinto explains, we can view the embedded method (2.21) as a Runge-Kutta method of $2s$ stages. Therefore, by lemma 2.4.4, a function $\rho(z)$ exists such that $y_n = [\rho(h\lambda)]^n$ for this method. We denote it by $\hat{R}(z)$. We can now impose the further condition on our coefficient in (2.21) that $\lim_{|z|\to\infty} \hat{R}(z) = \hat{R}(\infty) = 0$[¶]. This condition ensures the stability of the embedded method and is enforced through the following equation, as found in [17].

$$
\hat{R}(\infty) = 1 - \beta_n^T \mathbf{1} + \frac{(-1)^{s+1}}{s} \delta_n^T A^{-1} \mathbf{1} = 0
\tag{2.30}
$$

Solving this system of equations (2.22,2.23,2.30) for Lobatto IIIC coefficients (see section 2.6.2) with respect to $r = r_n$ with the condition that $\hat{p} = 3$ results in the following coefficients for the embedded method:

$$
\begin{array}{ll|ll}
\beta_1 = \frac{12r^3+14r^2+21r+9}{4r^3+4r^2+3r-3} & & \delta_1 = & \frac{12r^3+9r^2+3r}{4r^3+4r^2+3r-3} \\
\beta_2 = \frac{16r^3+8r^2-12r-12}{4r^3+4r^2+3r-3} & & \delta_2 = & \frac{8r^4-24r^3-36r^2-12r}{4r^3+4r^2+3r-3} \\
\beta_3 = 0 & & \delta_3 = & -\delta_1 - \delta_2
\end{array}
\tag{2.31}
$$

Note that $\beta_1, \beta_2, \delta_1, \delta_2, \delta_3$ all become undefined when $r = 0.5$. Plotting these coefficients we get the graph seen in Figure 2.1.
To avoid errors as $r$ gets close to 0.5 we redefine $r$ to be 0.4 if $r$ is in $[0.4, 0.6]$ and adjust $h_{n+1}$ accordingly before we take our step. This prevents any errors that could appear due to multiplying with large numbers and circumvents dividing by 0. Since this can only occur when the step size is currently decreasing, it just decreases it a bit extra. See the next section for more on the step size adjustment.

---

[¶]This is also beneficial to us since we do not fulfill Theorem 2.1 in [17] which implies that $\hat{R}(\infty)$ is bounded, hence by requiring this condition we ensure that it stays bounded

Figure 2.1: plotting the coefficient in (2.31)

Observer that the above error estimate does not work for the initial step, since we cannot calculate $r$. Thus, for the first step we must use a different method. Pinto et al. suggest solving the above equations with the added restriction $\delta_1 = \delta_2 = \delta_3 = 0$. However, since that would also result in an order 2 method we instead choose to use method (2.19) mentioned before. We now have an approximate error in each term as follows

$$err = \hat{y}_{n+1} - y_{n+1} \tag{2.32}$$

Note that both methods for $\hat{y}_{n+1}$ are of the form $y_n + \sum$. Similarly we have $y_{n+1} = y_n + \sum$. In fact for Lobatto IIIC with $s = 3$ we have $y_{n+1} = y_n + z_3$. We therefore change the error estimates to

$$err = h\gamma f(x_0, y_0) + \sum_{i=1}^{3} \hat{b}_i z_i - z_3 \tag{2.33}$$

for the first step and

$$err = a \left( \sum_{i=1}^{3} \delta_{n,i} z_{n-1,i} + \sum_{i=1}^{3} \beta_{n,i} z_{n,i} - z_3 \right) \tag{2.34}$$

for all other steps. The term $a$ is introduced by Pinto et al. and is a free parameter with the only restriction being $a > 0$. We chose a value of $a = 0.05$ which yielded good results, after doing some initial numerical experimentation with the problem discussed in section 4. However, a more thorough investigation into this parameter would be needed to optimize this error estimation and ensure that it works for a general problem.

## 2.5   Step Size Prediction

For this section we use $h_{new}$ synonymously to $h_{n+1}$ from the previous section on error estimation. The approximate error calculated in (2.32) is per term, we now need to combine these $n$ terms

into a single term. For this we use the norm introduced by Hairer in [14] :

$$||err|| = \sqrt{\frac{1}{n} \sum_{i=1}^{n} \left( \frac{err_i}{sc_i} \right)^2}$$ (2.35)

with $sc_i = Atol_i + \max(|y_{0,i}|, |y_{1,i}|) * Rtol_i$. The index $i$ refers to the $i$'th element in the corresponding vectors. $Atol$ and $Rtol$ are the desired absolute and relative tolerances per index in the resulting vector. The difference between this norm and a the two-norm is that this one weights the different indexes, giving one the possibility to be more precise in one index than the other.

Hairer then suggests the following step size control formula[†]

$$h_{new} = fac \cdot h_{old} \cdot ||err||^{-1/4}$$ (2.36)

The power of $-1/4$ stems from the order of the approximate solution $\hat{y}_{n+1}$ being 3, which means that the exact solution $\widetilde{y_n} = \hat{y}_n + \mathcal{O}(h^4)$. Note however, that in the first time step our approximate solution $\hat{y}_{n+1}$ is only of order 2 therefore in the first step instead of (2.36) we use

$$h_{new} = fac \cdot h_{old} \cdot ||err||^{-1/3}$$ (2.37)

The $fac$ term is dependent on the number of Newton iterations and is defined as

$$fac := 0.9 \times \frac{2k_{\max} + 1}{2k_{\max} + Newt}$$

with $k_{\max}$ being the maximum number of Newton iterations and $Newt$ being the actual number required until the stopping condition was reached. Hairer commented that they noticed best results setting $k_{\max}$ to around 7 or 10.

To not have to recompute the LU-decomposition on every step we shall keep the old step size $h_{old}$ if

$$c_1 h_{old} \leq h_{new} \leq c_2 h_{old}$$ (2.38)

Here Hairer suggests a value of 1.0 and 1.2 for $c_1$ and $c_2$ respectively. However, in our case, the step size seems to decrease continuously by a small fraction for the first 100 steps. Therefore we chose to set $c_1 = 0.97$. When the step size does not change and the Jacobian does not need to be recomputed then the LU-decomposition can be reused from the previous time step. This significantly reduces the number of LU-decompositions that need to be done.

## 2.6   The Classes of implicit Runge-Kutta Methods

We begin this section by giving a brief description of Radau and Lobatto method, and then examine the different Lobatto methods in more detail. As mentioned before the Runge-Kutta methods differ only in their choice of coefficients $a_{i,j}, b_i, c_i$. The choice of these coefficients stems

---

[†]Pito et al. suggest another method for step size control that could be applied as an extension However, is is beyond the scope of this work.

from the restrictions imposed on the method. As Hairer mentions one begins with the simplified assumptions

$$B(p): \sum_{i=1}^{s} b_i c_i^{q-1} = \frac{1}{q} \qquad\qquad q = 1,\ldots,p$$

$$C(\eta): \sum_{j=1}^{s} a_{i,j} c_j^{q-1} = \frac{c_i^q}{q} \qquad\qquad i = 1,\ldots,s,\ \ q = 1,\ldots,\eta$$

$$D(\xi): \sum_{i=1}^{s} b_i c_i^{q-1} a_{i,j} = \frac{b_j}{q}(1 - c_j^q) \qquad\qquad j = 1,\ldots,s,\ \ q = 1,\ldots,\xi$$

We now present, without a proof, a theorem connecting the simplified assumptions to the order of the Runge-Kutta method. A proof may be found in [8]

**Theorem 2.6.1.** *If the coefficients $b_i, c_i, a_{i,j}$ of a Runge-Kutta method satisfy $B(p), C(\eta), D(\xi)$ with $p \leq \eta + \xi + 1$ and $p \leq 2\eta + 2$, then the method is of order $p$.*

One can now use these conditions and aim for maximum order. Doing this results in the Gauss method of order $2s$. However, for this method the $c_i$ are not located at either of the end points. By imposing further restrictions on the $c_j$ terms we get the Radau and Lobatto methods. These methods have this name since they are based on the Radau and Lobatto quadrature formulas. The $c_j$'s where chosen according to the zeros of

$$\mathbf{I}: \frac{d^{s-1}}{dx^{s-1}}\left(x^s(x - 1^{s-1})\right) \qquad\qquad (2.39)$$

$$\mathbf{II}: \frac{d^{s-1}}{dx^{s-1}}\left(x^{s-1}(x - 1^s)\right) \qquad\qquad (2.40)$$

$$\mathbf{II}: \frac{d^{s-1}}{dx^{s-1}}\left(x^{s-1}(x - 1^{s-1})\right) \qquad\qquad (2.41)$$

Butcher called the Runge-Kutta methods with these $c_i$'s processes of type **I,II** and **III**, also known as Radau Left, Radau Right and Lobatto. The Radau left and right have those names since in Radau left the $c_0$ value is 0 and in Radau right the rightmost (last) $c_i$ value is 1. Which means that in Radau left we have a node at the left endpoint and in Radau right we have one at the right endpoint. Since these both have an extra restriction it is reasonable for their order to be less the optimum, which is known as the Gauss method and has the best possible order of $2s$. In fact they are both of order $2s - 1$. Lobatto on the other hand has an endpoint at both left and right side giving it one more restriction and an order of $2s - 2$ [14].

As mentioned before for this work we are interested specifically in the Lobatto IIIC method.

## 2.6.1   A, B & C

*The following section is based on the text by Laurent O. Jay.[16]*
The Lobatto method is not unique even with the choice of the $c_i$'s given by (2.41). However, since the $c_i$'s are defined the $b_i$'s they are also the same for all variants of Lobatto. That means the only coefficients that differ are the $a_{i,j}$'s.

- **A**:
  For the Lobatto IIIA method the coefficients are defined by $C(s)$, together with satisfying $D(s - 2)$ and the condition that $a_{s,j} = b_j$ and $a_{a,j} = 0$ for $j = 1,\ldots,s$.

- **B**:
  For the Lobatto IIIB method the coefficients are defined by $D(s)$ along with satisfying $C(s-2)$ and $a_{i,1} = b_1$ along with $a_{i,s} = 0$ for $i = 1, \ldots, s$.

- **C**:
  The Lobatto IIIC method is defined by $a_{i,1} = b_1$ and $C(s-1)$. They satisfy $D(s-1)$ and $a_{s,j} = b_j$ for $j = 1, \ldots, s$.

The difference between the A,B and C methods is in their stability. When talking about the stability of Runge Kutta methods one generally looks at weather they are $A$-Stable, $B$-Stable and/or $L$-stable. The different stability's look at different properties of the methods to analyse how well different numerical methods approximate the solutions to stiff ODEs. For example A-stability has to do with the stability function defined in (2.28). See [8] for more in-depth explanation on the different stability's. While all three are $A$-Stable neither A or B is $L$-stable or $B$-stable and hence neither is algebraically stable. The Lobatto IIIC method on the other hand is both $L$-stable and algebraically stable and hence also $B$-stable.

## 2.6.2  Lobatto IIIC

As mentioned before we are interested specifically in the Lobatto IIIC method. Additionally we focus on the three staged method, hence the coefficients we use are as follows [8, 14, 16]:

$$
\begin{array}{c|ccc}
0 & 1/6 & -1/3 & 1/6 \\
1/2 & 1/6 & 5/12 & -1/12 \\
1 & 1/6 & 2/3 & 1/6 \\
\hline
& 1/6 & 2/3 & 1/6
\end{array}
\tag{2.42}
$$

# Chapter 3

# Integration into Assimulo

Assimulo is a Python package that entails many ODE solvers written in different languages. It aims to provide a high-level interface for a variety of solvers. Thereby providing an easy way to compare different solvers for the same problem without having to redefine the problem in a different language to adjust for the different solvers [2]. We chose to implement the Lobatto IIIC solver as an extension to the Assimulo package for ease in the later comparison process of results and complexity to other solvers. In addition this solution allows for easier reuse of the code later on for other interested parties. We continue discussing some specifics of the implementation, for the full code consult the appendix.

## 3.1 The solver class

To define our own solver in Assimulo all we are required to do is to define our own class as a child of `Explicit_ODE`. This class should overload the `integrate` method which is called when the solver is asked to solve a problem. We have split our implementation into two classes, the `RKS3CBase` and `Lobatto4ODE` with the latter being a child of the former. This was done so that we could easily implement different Runge-Kutta methods later.

### 3.1.1 Integrate method

The `integrate` method is defined as `def integrate(self,t,y,tf,opts)` with `t` being the initial time $x_0$, `y` being the initial value $y_0$, `tf` being the final time and `opts` being the options. The method is then expected to return the `ID_PY_OK` object along with a list of time stamps and a corresponding list or approximations for each time stamp, `tres` and `yres` respectively in our case. Note that we did not use the options in our implementation. These options could be used to specify that we only wish to receive outputs at certain time points but wish to have a smaller step size, i.e. only store some points. Doing this would help reducing the memory requirements of the entire Runge-Kutta method. The `integrate` method itself is just a loop, taking a step and then storing the results of the new $x_n$ and $y_n$ in the `tres` and `yres` lists respectively. The `step` method is just a helper method which calls the `_newton` method and passes the results back to the `integrate`. The reason that we have the `step` method at all is that the `_newton` sometimes calls itself recursively when the step fails.

### 3.1.2   Newton Iteration method

The `_netwon` method is the main part of the code. This implements the entire Newton iteration described in section 2.2 together with all the simplifications from section 2.3. To perform the newton iteration described in equation (2.14) we require the left hand side (LHS) and right hand side (RHS) of equation (2.14a) so that we can solve for $\Delta W^k$. The computation of these sides is done in the methods `_compLHS` and `_compRHS` respectively.

The `_compLHS` method constructs the blocks for the block matrix in (2.17). This requires only the Jacobian $\mathscr{J}$ of the problem, which can be passed as a sparse or dense matrix. This method is only called once per newton iteration and only if the step size $h$, called `h` in the code, has been adjusted. In addition once `_compLHS` is called the result gets passed into the LU-decomposition method. The result of the LU-decomposition is then stored as an instanced variable to be reused in later newton iterations.

The `_compRHS` computes $-h^{-1}(\Lambda \otimes I)W^k + (T^{-1} \otimes I)F(Z^k)$, see equation (2.14a), which only depends on the current $Z^k, W^k$ and the function $f$ of the ordinary differential equation. This method is called once per newton step, and has three evaluations of the function $f$. It then returns a single numpy vector of size $3n$

The `_netwon` method also evaluates the error after the Newton iteration terminates with `err_mag = self._normed_aprox_err(zk,y0,y1,n)`. Further it adjusts the step size or restarts the step if the error was to large. The method `_normed_aprox_err` is defined within the Lobatto class and is error approximation method.

# Chapter 4

# Numerical Results

## 4.1 Rotating Pulse Problem

To test our implementation we used the rotating pulse problem discussed in Section 5.3 in [5].

$$\frac{\delta C}{\delta t} - \nabla \cdot \left\{ \begin{pmatrix} -4y \\ 4x \end{pmatrix} C - D \nabla C \right\} = 0 \tag{4.1}$$

with the Dirichlet boundary and initial conditions from the analytical solution

$$C(x, y, t) = \frac{2\sigma^2}{2\sigma^2 + 4Dt} \exp\left( -\frac{(\bar{x} - x_c)^2 + (\bar{y} - y_c)^2}{2\sigma^2 + 4Dt} \right) \tag{4.2}$$

with $\bar{x} = x\cos(4t) + y\sin(4t)$, $\bar{y} = -x\sin(4t) + y\cos(4t)$. We shift the domain to $\Omega = (0, 1)^2$ instead of $(-0.5, 0.5)^2$ and use the parameters $D = 10^{-3}, x_c = 0.25, y_c = 0$ and $2\sigma^2 = 0.004$. The time interval for the simulation is $[0, 0.25]$. This problem describes a rotating heat pulse which cools down over time as can be seen in Figure 4.1 where we show 8 instances in time extending the simulation time to see the effect more clearly. The benefit of using this specific problem is that an analytical solution exists. On that account we are able to directly measure the error between approximation and analytical solution.

We used an implementation of the problem in the DUNE framework [6, 3, 4, 9, 10, 1, 12]. Then we pass the RHS, $f$, of our ODE, along with the jacobian $\mathcal{J}$ from DUNE to Assimulo. To define a ODE problem in Assimulo is as simple as calling `Explicit_Problem(rhs,y0)`. Which returns a model object. Here `rhs` is the function, $f$, in our code and should have inputs `rhs(t,y)`. While the input `y0` is the initial condition, $y_0 = f(x_0)$, of our ODE. Unless an addition parameter for $x_0$ is passed in it is assumed to be zero. To then find the solution at $x_N \neq x_0$ one first defines a solver, in our case `Lobatto4ODE(model)` which returns a simulator. Here `model` is the model object returned from the Explicit Problem. To then simulate this simulator until $x_N$ we simply call `.simulate(end)` where `end` is $x_N$. This returns a list of time points and a corresponding list of solutions.

We run the code* with two different spatial resolutions, the first discretization discretizes the space into a grid of $20 \times 20$ points and the second having $2 \cdot 20 \times 2 \cdot 20$ points, referred to as *resolution 1* and *resolution 2* respectively from here on. As a point of comparison we will compare the error from our method to the error in the RADAU5 method implemented by Hairer in [14],

---

*See appendix for the code used to setup and solve the problem.

which has already been implemented into the Assimulo package[2]. Since our Lobatto method is of order 4 while the already existing Radau method is of order 5 we expect the Radau method to outperform our method. For a more fair comparison we will also compare the to RungeKutta34 method in Assimulo. This method is of order 4 and as such should behave similar to our method.



(a) t=0.25     (b) t=0.5     (c) t=0.75     (d) t=1.00

(e) t=1.25     (f) t=1.5     (g) t=1.75     (h) t=2.00

Figure 4.1: A sequence of results from the numerical experiment at different points in time showing the movement of the pulse.

## 4.2 Results

| | Lobatto | | Radau5 | | RungeKutta34 | |
|---|---|---|---|---|---|---|
| Resolution | 1 | 2 | 1 | 2 | 1 | 2 |
| Computation time (sec) | 7.98 | 159 | 11.36 | 403 | 0.69 | 17.61 |
| Number of function calls | 70 | 76 | 62 | 58 | 380 | 980 |
| Number of Jacobi evaluations | 1 | 1 | 1 | 1 | - | - |
| Number of LU-decompositions | 6 | 9 | 4 | 4 | - | - |
| Number of steps | 18 | 18 | 14 | 13 | 76 | 196 |
| $L^2$ Error | 0.00225 | 0.000244 | 0.00226 | 0.000254 | 0.00226 | 0.000251 |
| Average number of function calls per steps | 3.888 | 4.222 | 4.429 | 4.462 | 5.000 | 5.000 |

Table 4.1: A comparison between the Lobatto IIIC, Radau IIA and RungeKutta34 method when simulating the rotating pulse problem as described above with a relative and absolute tolerance of $10^{-5}$



Figure 4.2: Showing the relationship between computation time in seconds and the achieved accuracy of our Lobatto IIIC method. Observe that both axes are in logarithmic scale

| Resolution | 1 | 2 |
|:---:|:---:|:---:|
| Computation time (sec) | 3.68 | 43.7 |
| Number of function calls | 60 | 60 |
| Number of Jacobi evaluations | 1 | 1 |
| Number of LU-decompositions | 2 | 2 |
| Number of steps | 16 | 16 |
| $L^2$ Error | 0.00226 | 0.000250 |
| Average number of function calls per steps | 3.75 | 3.75 |

Table 4.2: Comparison between resolutions when running our implementation of the Lobatto IIIC method with a constant step size for the rotating pulse problem as described above

## 4.3   Comparison of Methods

Table 4.1 shows the difference in the stats gathered while simulating the rotating pulse problem between our Lobatto implementation, Hairers Radau5 and the RungeKutta34 method in Assimulo. We continue by discussing their differences in the following sections.

### 4.3.1   $L^2$ Error

The $L^2$ error of all the methods does not go down to zero, this is because we are taking the norm between our result and a discritization of the exact solution, this discritization introduces a small error. This is also the reason why the error for all methods is nearly the same.

### 4.3.2   Computational Time

The computational time is the time that our computer clock measures between starting the simulation and finishing. The difference in computational time between the different methods is mainly due to the differences in the implementation of the methods. Making this a bad metric to judge the different methods by. Our Lobatto code for example is written completely in Python while the Radau5 code is written in FORTRAN and we only call it in from Python. Since Python is an interpreted language instead of a compiled one it is slightly slower and therefore will take more time per step. In addition the Radau5 code is unable to use the sparse matrix representation used for the Jacobian and hence for each Jacobian evaluation the matrix must be converted to dense representation first. This might be why for resolution 2 the Radau5 takes longer to compute compared to our Lobatto method. Similarly taking an LU decomposition of a large dense matrices takes longer in comparison to sparse matrices if the matrix is large enough and is filled with significantly many zeros.

We have further investigated the computation time of our Lobatto IIIC method compared to the achieved error, see figure 4.2. To generate this graph we changed the resolution of the discretization. The finer the grid the lower the error however at the cost of computation time.

### 4.3.3   Number of Function Calls

The number of function calls is directly related to how many steps were taken and how many of those steps failed. As can be seen in all three methods the more steps one needs the more function calls are required. However the Radau5 method has at least one extra function call per step through its error evaluation while the Lobatto method does not. Hence, if we look at the average number of function calls per step we see that our Lobatto method has fewer. Since we only solve for $t = 0.25$ very few steps need to be taken however as the end time increases so dose the number of steps, and as this problem is linear the average number of steps in Lobatto will go to 3 while in Radau it will go to 4.

### 4.3.4   Number of Jacobian Evaluations

Due to the structure of this problem a single Jacobian evaluation is enough and both our Lobatto and the Radau5 method utilize this.[†]

---

[†]The RungeKutta34 method probably also uses this considering how fast it is.

### 4.3.5 Number of LU-decompositions

Our Lobatto method has a few more LU-decompositions than the Radau5 method. The decomposition is recalculated whenever the step size was adjusted. This means that the Radau5 method was slightly faster at determining the required step size[‡].

### 4.3.6 Number of Steps

The number of steps is directly linked to the step size. Since the Lobatto method is not as good as the Radau method in terms of order, more steps will need to be taken to achieve the same accuracy in the results. Similarly the RungeKutta34 method is, like our Lobatto, order 4. Hence we would expect to see a similar amount of steps. The spatial discritization we used is optimized for the Lobatto points, which is probably why the Lobatto Method requires slightly fewer steps than RungeKutta34.

### 4.3.7 Overall

Over all our Lobatto method has performed as expected. Being slightly better than the RungeKutta34 but still not as good as the Radau5 method in all metrics other than computation time. However, computation time, as explained is not a good metric since it would require all methods to be implemented in the same programming language with the same amount of optimization. Note further for both our implementation of the Lobatto method and the Radau5 method from Hairer [14] neither of the two significantly increase as the resolution increases, however, the RungeKutta34 method more than doubles the amount of steps it takes.

## 4.4 Constant step

The interest of this project was sparked due to the result in [19] stating that the Lobatto IIIC method is equivalent to the space time DG approximation. When solving for a space time discritization, the step size is the dicritization in time and is kept constant. We therefore run the before mentioned problem again with a fixed step size of $h = \frac{1}{60}$, and get the results shown in table 4.2. This table shows, that when solving linear problems, manually determining a constant step size can be beneficial. The fact that a lower number of steps still gives an accurate result is probably due to an suboptimal choice of our free parameter $a$ in equation (2.34)

---

[‡]An idea exists to only recalculate the LU-decomposition when Newton diverges, this was not done here however

# Chapter 5

# Conclusion

In this work we have outlined the implementation of the Lobatto IIIC Runge-Kutta method as described by Hairer in [14]. We have then done a more in-depth discussion about the different existing error estimates for the local truncated error, and chosen to adapt a method from Radau to Lobatto that Pinto et al. describes in [17].

Comparing our method to the Radau5 method from Hairer the only improvement we see is a lower average amount of function calls per step. However as discussed the Radau5 method has a higher order and hence is expected to outperform our Lobatto IIIC implementation. When comparing our method to the RungeKutta34 method instead, we see that we have significantly fewer steps and hence significantly fewer function calls. Be that as it may, we have only conducted one numerical experiment comparing the methods and further experiments must be conducted with varying ODEs for a thorough comparison.

Further improvements could still be made to our method, such as thorough investigating into a good choice for the free parameter $a$ in equation (2.34) and the multitude of extensions mentioned throughout this work which were excluded due to the limited scope of this study.

# Chapter 6

# Bibliography

[1] M. Alkämper et al. "The DUNE-ALUGrid Module." In: *Archive of Numerical Software* 4.1 (2016), pp. 1–28. ISSN: 2197-8263. DOI: 10.11588/ans.2016.1.23252. URL: http://dx.doi.org/10.11588/ans.2016.1.23252.

[2] Christian Andersson, Claus Führer, and Johan Åkesson. "Assimulo: A unified framework for {ODE} solvers". In: *Mathematics and Computers in Simulation* 116.0 (2015), pp. 26–43. ISSN: 0378-4754. DOI: http://dx.doi.org/10.1016/j.matcom.2015.04.007.

[3] P. Bastian et al. "A Generic Grid Interface for Parallel and Adaptive Scientific Computing. Part I: Abstract Framework". In: *Computing* 82.2–3 (2008), pp. 103–119. DOI: 10.1007/s00607-008-0003-x. URL: http://dx.doi.org/10.1007/s00607-008-0003-x.

[4] P. Bastian et al. "A Generic Grid Interface for Parallel and Adaptive Scientific Computing. Part II: Implementation and Tests in DUNE". In: *Computing* 82.2–3 (2008), pp. 121–138. DOI: 10.1007/s00607-008-0004-9. URL: http://dx.doi.org/10.1007/s00607-008-0004-9.

[5] Peter Bastian. "Higher Order Discontinuous Galerkin Methods for Flow and Transport in Porous Media". In: *Challenges in Scientific Computing- CISC 2002*. Ed. by Eberhard Bänsh. Springer, 2002, pp. 1–22. ISBN: 978-3-642-19014-8.

[6] Peter Bastian et al. "The Dune framework: Basic concepts and recent developments". In: *Computers  Mathematics with Applications* 81 (2021). Development and Application of Open-source Software for Problems with Numerical PDEs, pp. 75–112. ISSN: 0898-1221. DOI: https://doi.org/10.1016/j.camwa.2020.06.007. URL: https://www.sciencedirect.com/science/article/pii/S089812212030256X.

[7] Phillip Birken. *Newton's Method 2*. May 2020.

[8] J. C. Butcher. *The numerical analysis of ordinary differential equations*. John Wile  Sons, 1987. ISBN: 0 471 91046 5.

[9] A. Dedner et al. "A Generic Interface for Parallel and Adaptive Scientific Computing: Abstraction Principles and the DUNE-FEM Module". In: *Computing* 90.3–4 (2010), pp. 165–196. DOI: 10.1007/s00607-010-0110-3. URL: http://dx.doi.org/10.1007/s00607-010-0110-3.

[10] A. Dedner et al. "The DUNE-FEM-DG module". In: *Archive of Numerical Software* 5.1 (2017), pp. 21–61. ISSN: 2197-8263. DOI: 10.11588/ans.2017.1.28602. URL: http://journals.ub.uni-heidelberg.de/index.php/ans/article/view/28602.

[11]  Andreas Dedner, Robert Kloefkorn, and Martin Nolte. *Python Bindings for the DUNE-FEM module*. Version v2.7.0. Mar. 2020. DOI: 10.5281/zenodo.3706994. URL: https://doi.org/10.5281/zenodo.3706994.

[12]  Andreas Dedner et al. *Python Framework for HP Adaptive Discontinuous Galerkin Method for Two Phase Flow in Porous Media*. 2018. arXiv: 1805.00290 [cs.CE].

[13]  Brian C. Fabien and P. Frye Jason. "The solution of the inplicit differential equation via a lobatto IIIC Runge-Kutta method". In: *MATHMOD09*. 2009, pp. 2708–2711. ISBN: 978-3-901608-35-3. URL: https://www.argesim.org/fileadmin/user_upload_argesim/ARGESIM_Publications_OA/MATHMOD_Publications_OA/MATHMOD_2009_AR34_35/full_papers/303.pdf.

[14]  Ernst Hairer and Gerhard Wanner. *Solving Ordinary Differential Equations II*. 2nd ed. Vol. 2. Springer, 2010, pp. 135–144.

[15]  Arieh Iserles. *A First Course in the Numerical Analysis of Differential Equations*. 2nd ed. Cambridge University Press, 2009. ISBN: 978-0-511-50637-6.

[16]  O. Jay Laurent. *Lobatto methods*. URL: https://homepage.math.uiowa.edu/~ljay/publications.dir/Lobatto.pdf.

[17]  Severiano González Pinto, Domingo Hernández Abreu, and Juan Ignacio Montijano. "Variable Step-Size Control Based on Two-Steps for Radau IIA Methods". In: *ACM Transactions on Mathematical Software* 46.4 (2020). DOI: 10.1145/3408892.

[18]  Endre Süli and David Francis Mayers. *An introduction to numerical analysis*. Cambridge University Press, 2014. ISBN: 957-0-521-810267-5.

[19]  Lea Miko Versbach. *Multigrid Preconditioners for the Discontinuous Galerkin Spectral Element Method*. 2020. ISBN: 978-91-7895-593-0.

[20]  Wikipedia contributors. *List of Runge–Kutta methods — Wikipedia, The Free Encyclopedia*. [Online; accessed 21-May-2021]. 2021. URL: https://en.wikipedia.org/w/index.php?title=List_of_Runge%E2%80%93Kutta_methods&oldid=1019271082.

# Chapter 7

# Appendix

## 7.1 Lobatto IIIC Implementation

```python
1  import sys #for min float number
2  import numpy as np
3  import scipy.sparse as sp
4  import scipy.sparse.linalg as alg
5  from scipy.sparse.linalg import splu
6  from abc import ABC, abstractclassmethod
7  from scipy.linalg import lu_factor,lu_solve
8
9  from assimulo.explicit_ode import Explicit_ODE
10 from assimulo.ode import ID_PY_OK,NORMAL
11
12
13
14 class RadauError(Exception):
15     def __init__(self,expr):
16         self.expr = expr
17     def __str__(self):
18         return self.expr
19
20 class RKS3CBase(ABC,Explicit_ODE): #Runge-Kutta S=3 Complex eigen pair Base class
21
22
23
24     rtol=1.e-4 #realative tolerance
25     atol=1.e-4 #absolute tolerance
26     tol =1.e-4 #tollerance for newton itteration
27     kappa = 1e-1 #The kappa value used for newton stopping in eq (15)
28     maxit=10 #maximum number of newton itterations
29     maxsteps=10000000 #maximum number of steps
30     r=1 #r is the ratio between the previouse step size and the current
31
32     #used for stoping after single newton itt
33     Uround = sys.float_info.min
34
35     #the bounds to decide weather or not to adjust the step size eq (41)
36     c1 = 0.95
37     c2 = 1.2
```

```
38
39
40     def __init__(self,problem):
41         """
42         Initiates the solver.
43
44             Parameters::
45
46                 problem
47                         - The problem to be solved. Should be an instance
48                           of the 'Explicit_Problem' class.
49         """
50         Explicit_ODE.__init__(self, problem) #Calls the base class
51
52         #Solver options
53         self.options['h'] = problem.h
54         self.h = problem.h
55         self._sparceJ = True
56         self._dynamic_step = True
57
58         #Statistics
59         self.statistics["nsteps"] = 0
60         self.statistics["nfcns"] = 0
61         self.num_jac_evals = 0
62         self.num_lu_decomps = 0
63
64         #method parameter:
65         self.y_dim = len(self.problem.y0)
66         self.RTol = np.array([self.rtol]*self.y_dim)
67         self.ATol = np.array([self.atol]*self.y_dim)
68         self._3ATol = np.hstack((self.ATol,self.ATol,self.ATol))
69         self._3RTol = np.hstack((self.RTol,self.RTol,self.RTol))
70
71         #internal flags
72         self._first = True
73         self._recomp_jac = True
74         self._recomp_LHS = True
75         self._newt_fail_flag = False
76
77         #internal values
78         self.J = None
79         self._lu = None
80
81
82     #set and get for sparce Jacobian, need to redifine the lu decompose and
83     #lu solve methods here
84     def _set_sparce(self,sp):
85         self.options["sparce"] = bool(sp)
86         if sp:
87             def lu_decomp_s(LHS):
88                 lu1 = splu(LHS[0])
89                 lu2 = splu(LHS[1])
90                 return (lu1,lu2)
91             self._LUDecomp = lu_decomp_s
92
93             def lu_solve_s(lu,RHS):
```

```python
 94                    n = RHS.shape[0]//3
 95                    dwk = np.zeros_like(RHS)
 96                    dwk[:n] = lu[0].solve(RHS[:n])
 97                    dwk[n:] = lu[1].solve(RHS[n:])
 98                    return dwk
 99                self._LUSolve = lu_solve_s
100            else:
101                def lu_decomp_d(LHS):
102                    lu1 = lu_factor(LHS[0])
103                    lu2 = lu_factor(LHS[1])
104                    return (lu1,lu2)
105                self._LUDecomp = lu_decomp_d
106
107                def lu_solve_d(lu,RHS):
108                    n = RHS.shape[0]//3
109                    dwk = np.zeros_like(RHS)
110                    dwk[:n] = lu_solve(lu[0],RHS[:n])
111                    dwk[n:] = lu_solve(lu[1],RHS[n:])
112                    return dwk
113                self._LUSolve = lu_solve_d
114        def _get_sparce(self):
115            return self.options["sparce"]
116        _sparceJ = property(_get_sparce,_set_sparce)
117
118        #set and get for h parameter, when h changes we also want to recompute the
119        #LHS
120        def _set_h(self,h):
121            self.options["h"] = float(h)
122        def _get_h(self):
123            return self.options["h"]
124        h=property(_get_h,_set_h)
125
126
127        def integrate(self,t,y,tf,opts):
128            """
129            _integrates (t,y) values until t > tf
130            """
131            self.h = min(self.h,abs(tf-t))
132
133            tres = []
134            yres = []
135
136            z0 = np.zeros((self.y_dim*3))
137
138            #lists for storeing results
139            for i in range(self.maxsteps):
140                if t>= tf:
141                    break
142
143                t,y = self._step(t,y,z0)
144
145                self._first = False
146
147                self.t = t #used for Radau error approximation
148
149                #store the results in array to be returned when done
```

29

```
150                tres.append(t)
151                yres.append(y.copy())
152
153
154            if self.h > np.abs(tf-t):
155                h_old = self.h / self.r
156                self.h = np.abs(tf-t)
157                self.r = self.h/h_old
158                self._recomp_LHS=True
159
160            if abs(self.r-0.5)<0.1:
161                #if we are to close to the explosion point, change value
162                h_old = self.h/self.r
163                self.r = 0.35
164                self.h = self.r * h_old
165                self._recomp_LHS=True
166
167            if abs(self.h) < self.Uround:
168                break
169
170        else:
171            raise RadauError(f"final time not reached within maximum number of steps ({self.maxsteps})")
172        return ID_PY_OK, tres, yres
173
174
175    def _step(self,t,y,z0):
176        #update statistics
177        self.statistics["nsteps"] += 1
178        #recomp jacobian if the flag is set
179        if self._recomp_jac:
180            self.J = self._jacobian(t,y)
181            self._recomp_jac = False
182            self._recomp_LHS = True
183
184        #newton itteration
185        y,h = self._newton(self.J,y.copy(),t,z0)
186        return t + h, y
187
188    def _newton(self,J,y0,t,z0=None):
189
190        #handy dimension variable
191        n  = self.y_dim
192
193        #setup newton starting value
194        zk = z0 if type(z0) is np.ndarray else np.zeros(3*n)
195        wk = self._zTOw(zk)
196
197        #newton stopping variables
198        ndzk = None #Norm of delta Zk
199        theta_k = None #Theta_k
200
201        #setup LHS and lu decompose if needed
202        if self._recomp_LHS or self._lu == None:
203            LHS = self._compLHS(J)
204            self._lu = self._LUDecomp(LHS)
205            self.num_lu_decomps += 1
```

30

```python
206              self._recomp_LHS = False
207
208         for i in range(1,self.maxit):
209             #compute the RHS
210             RHS = self._compRHS(wk,zk,y0)
211
212             #compute the delta W_k
213             dwk = self._LUSolve(self._lu,RHS) #eq (17a)
214
215             #convert delta W_ktodelta Z_k
216             dzk = self._wTOz(dwk)
217
218             #increment wk and zk with the corresponding delta values
219             wk += dwk #eq (17b)
220             zk += dzk
221
222             #---------newton stopping condition-----------
223
224             #for sci we require the max of z_k and z_{k-1} termwise
225             diff = dzk.copy()
226             diff[diff>0] = 0 #wherever the value has increased we want to keep
227             sci = self._3ATol + (zk-diff)*self._3RTol
228
229             #norm of delta zk
230             ndzk = self._norm(dzk,sci)
231
232             #declare the eta_k variable so it is outside the if scope
233             eta_k = None
234             if i == 1 :#first step in newton
235                 if '_eta' in self.__dict__:
236                 #one step stopping
237                     eta_k = max(self._eta,self.Uround)**.8
238                 else:
239                 #first time step, needs at least 2 newton
240                     ndzkm1 = ndzk
241                     continue
242             else:
243                 #compute theta_k and eta_k
244                 theta_k = ndzk/ndzkm1 # eq (14)
245                 eta_k   = theta_k/(1-theta_k) #eq (15)
246
247                 #check if we are diverging
248                 if theta_k > 1:
249
250                     if not self._dynamic_step:
251                         raise RadauError(f'newton itteratio diverged')
252
253                     #adjust step size
254                     self.h /= 2
255                     self.r /= 2
256
257                     #if r is to close to 0.5 adjust it down further,
258                     #this is due to explosion of values in error estimation
259                     #TODO: move this into the set h function.
260                     if (self.r-0.5)<0.1:
261                         self.h *= 0.35 /self.r
```

```
262                         self.r = 0.35
263
264                     #recompute the jacobian
265                     self.J = self._jacobian(t,y0)
266                     self._recomp_LHS = True
267
268                     #re-run newton
269                     return self._newton(self.J,y0,t,z0)
270
271             #test if newton has converged, eq (14)
272             if eta_k * ndzk <= self.kappa * self.tol:
273                 #save old eta for singe step stopping
274                 self._eta = eta_k
275
276                 #new y value
277                 y1 = self._zTOy(y0,zk)
278                 #the h value that was used
279                 h = self.h
280
281                 #------determine approximate error---------
282                 if self._dynamic_step:
283                     #approximate error magnitude
284                     err_mag = self._normed_aprox_err(zk,y0,y1,n)
285
286                     #if error magnitude > 1 last h will not be changed since
287                     #we will rerun this step.
288                     h_old = self.h/ ( self.r if err_mag > 1 else 1)
289
290                     #compute new step size
291                     h_new = self._adjust_h(err_mag,step_accept = err_mag < 1, itt = i)
292                     #test if we should adjust step eq (41):
293                     if not (self.c1*self.h <= h_new and h_new <= self.c2 * self.h):
294                         #set flag to recompute the LHS in the next itteration
295                         self._recomp_LHS = True
296                         #calculate the new ratio
297                         r = h_new/h_old
298
299                         #special condition for the Lobatto to avoid coefficient
300                         #explosion when computeing the error magnitude
301                         if abs(r-0.5)<0.1:
302                             #if we are to close to the explosion point, change value
303                             r = 0.35
304                             h_new = r * h_old
305                         #save adjusted values
306                         self.r = r
307                         self.h = h_new
308                     else:
309                         #if we dont adjust step size
310                         h_new = self.h
311                         self.r = 1.0
312
313
314                 if  err_mag > 1:
315                     #approximate error is to large, restart with smaller step size
316                     return self._newton(J,y0,t,z0)
317
```

```
318                      #store z^(k-1) for error estimation in the next step.
319                      self.zkm1 = zk.copy()
320
321                  #return result and the h used.
322                  return y1, h
323
324              #update the norm for newton stopping condition
325              ndzkm1 = ndzk
326
327          else:
328              #if we do not converge within number of newton itteration steps
329              if not self._newt_fail_flag:
330                  self._newt_fail_flag = True
331                  #rerun with new Jacobian if not done yet
332                  self.J = self._jacobian(t,y0)
333                  self._recomp_LHS = True
334                  #re-run newton
335                  return self._newton(self.J,y0,t,z0)
336              #otherwise raise
337              raise RadauError(f'newton itteratio did not converge after {self.maxit} itterations')
338
339      @abstractclassmethod
340      def _normed_aprox_err(self,zk,y0,y1,n):
341          #abstaract method that each version of this code must impliment seperately
342          return np.Inf
343
344      @abstractclassmethod
345      def _adjust_h(self,err_mag,TOL = 1e-5,step_accept = True,itt = 0):
346          return np.Inf
347
348      def _norm(self,vec,sc):
349          #calculates the norm of a vector according to eq (38)
350          return np.sqrt(
351                  (1/vec.shape[0])*
352                  np.sum(
353                      (vec/sc)**2
354                  )
355              )
356
357
358
359      def _compRHS(self,wk,zk,y):
360          #computes the RHS of eq (17a)
361
362          #dimenson of system
363          n = self.y_dim
364          #initial term:
365          term1 = -np.array(n*[self.gamma/self.h] + \
366                          2*n*[(self.alpha)/self.h]) \
367                      * wk
368          term1[n:2*n] += self.betta/self.h * wk[2*n: ]
369          term1[2*n: ] -= self.betta/self.h * wk[n:2*n]
370
371          #function term
372          fcn = self.problem.rhs
373          f = np.array([fcn(self.t+self.c[0]*self.h,y + zk[:n]),
```

```
374                fcn(self.t+self.c[1]*self.h,y + zk[n:2*n]),
375                fcn(self.t+self.c[2]*self.h,y + zk[2*n:])])
376            prod = (self.TI @ f).flatten()
377
378            self.statistics["nfcns"] += 3
379
380            return term1 + prod
381
382        def _compLHS(self,J):
383            #construction of the LHS of eq (17a)
384            #dimension of system
385            n = self.y_dim
386
387            #construct of blocks according to eq (20)
388            e1 = np.diag(n*[self.gamma/self.h])-J
389            e2 = np.zeros((2*n,2*n))
390            np.fill_diagonal(e2,self.alpha/self.h)
391            e2[:n,:n]  -= J
392            e2[n:,n:]   = e2[:n,:n]
393            np.fill_diagonal(e2[ :n,n: ],-self.betta/self.h)
394            np.fill_diagonal(e2[n: , :n], self.betta/self.h)
395
396            #if sparce then convert to sparce format
397            if self._sparceJ:
398                e1 = sp.csc_matrix(e1)
399                e2 = sp.csc_matrix(e2)
400            return (e1,e2)
401
402        def _jacobian(self,t,y):
403            """
404            calculates the jacobian
405            """
406            if self.problem_info['jac_fcn']:
407                self.num_jac_evals += 1
408                J = self.problem.jac(t,y)
409                self._recomp_jac = False
410                return J
411            else:
412                raise RadauError('The current implimentation cannot calculate an approximate jacobian yet,' +
413                        'hence a jacobian function must be given in the problem')
414
415        def _zTOw(self,z):
416            #TODO: this might one returns return directly...
417            #convert from z to w
418            w = (self.TI@ z.reshape((3,-1))).flatten()
419            return w
420
421        def _wTOz(self,w):
422            #convert from w to z
423            z = (self.T @ w.reshape((3,-1))).flatten()
424            return z
425
426        def _zTOy(self,y,z):
427            #convert from z to y
428            y = y + self.d @ z.reshape((3,-1))
429            return y
```

```python
430
431     def print_statistics(self, verbose=NORMAL):
432         """
433         Prints the run-time statistics for the problem.
434         """
435         Explicit_ODE.print_statistics(self, verbose) #Calls the base class
436         self.log_message(' Number of J evaluations : ' + str(self.num_jac_evals),    verbose)
437         self.log_message(' Number of LU decompos   : ' + str(self.num_lu_decomps),   verbose)
438         self.log_message('\n\n\n',verbose)
439
440
441 class Lobatto4ODE(RKS3CBase):
442
443     #-------- Method constants -----------
444     # T matrix
445     T = np.array([
446         [0.4072639531732107,-0.44308062047843144,0.31680119776208315],
447         [0.18547209365357897,0.1305271017756723,-0.38187534481524843],
448         [0.8942796961362183,0.735153359223255,0.0]
449         ])
450     # T inverse matrix
451     TI= np.array([
452         [1.0326372242409414,0.8566688421623923,0.47027336073401904],
453         [-1.2561549117980622,-1.0420976007887581,0.7881948366175593],
454         [0.07217833861373389,-2.558776912784013,0.49781525580043673],
455         ])
456
457     #eigen values gamma, alpha and betta
458     gamma = 2.6258168189584676
459     alpha = 1.6870915905207662
460     betta = -2.508731754924879
461
462     # c,d and b (see paper for reference)
463     c = np.array([0.,1/2.,1.])
464     d = np.array([0,0,1])
465     b = np.array([1/6,2/3,1/6])
466
467     #e used for one step error approximation
468     #e = np.array([-.5,2,-.5])-b
469     e = np.array([-4.0,0,-1]) #this is e A^-1
470
471     #free parameter a used in error estimation
472     a = 5
473
474     def __init__(self,problem):
475         RKS3CBase.__init__(self,problem)
476
477         #list for storeing error compairison if flag is set
478         self.err_comp = []
479
480     #overwrite z to y conversion to be more efficient
481     def _zToy(self,y,z):
482         return y + z[2*self.y_dim:]
483
484
485     #Method to calculate the coefficients for the error estimation
```

```python
486     def coef(self,r):
487         if abs(r-0.5)<0.09:
488             raise RadauError(f'Some internal error; the relative distance in step size is to close to 0.5,
489
490         #TODO: Optimize this method
491         #-devisor is always the same
492         #-delta3 = -delta1 -delta2
493         #-precalculate r^3 and r^2 (this might make it more efficient?)
494         r2 = r**2
495         r3 = r*r2
496         r4 = r*r3
497         div = 1/(4*r3 + 4*r2 + 3*r - 3)
498         beta1 = (12*r3 + 14*r2 + 21*r + 9)*div
499         beta2 = (16*r3 + 8*r2 - 12*r - 12)*div
500         beta3 = 0
501         delta1=(12*r3 + 9*r2 + 3*r)*div
502         delta2=(8*r4 - 24*r3 - 36*r2 - 12*r)*div
503         #delta3=(-8*r4 + 12*r3 + 27*r2 + 9*r)*div
504         return (beta1,beta2,delta1,delta2,-delta1-delta2)
505
506
507     def _normed_aprox_err(self,zk,y0,y1,n):
508         #calculate scaleing of error in the norm from hairer
509         sc = self.ATol + np.maximum(np.abs(y0),np.abs(y1))*self.RTol
510
511         #if we are in the first step
512         if self._first or self._do_comp_err:
513             #do one step error estimation
514             zpart = self.e[0] * zk[:n] + self.e[1] * zk[n:2*n] + self.e[2] * zk[2*n:]
515             ydiff = self.h*self.problem.rhs(self.t,y0)/self.gamma + zpart
516             self.statistics["nfcns"] += 1
517
518             #normalize error vector
519             mag = self._norm(ydiff,sc)
520
521             #if we want to compair the error to the two step error
522             if (not self._first):
523                 beta1,beta2,delta1,delta2,delta3 = self.coef(self.r)
524                 err_vec_tmp = self.a * (zk[2*n:] - (delta1*self.zkm1[:n] \
525                     + delta2*self.zkm1[n:2*n] + delta3*self.zkm1[2*n:] \
526                     + beta1 * zk[:n] + beta2 * zk[n:2*n]))
527                 err = self._norm(err_vec_tmp,sc)
528                 err_old = mag
529                 #return err
530
531                 self.err_comp.append((err ,mag,self.t,self.h))
532                 if mag > 1:
533                     self.statistics["nfcns"] += 1
534                     mag = self._norm(self.h * self.problem.rhs(self.t,y0)/self.gamma + zpart,sc)
535
536
537             else:
538                 err = mag
539                 #if the step fails do more accurate error estimation
540                 if err > 1:
541                     self.statistics["nfcns"] += 1
```

```
542                        err = self._norm(self.h * self.problem.rhs(self.t,y0)/self.gamma + zpart,sc)
543                return err
544
545            else:
546                #do two step error estimation
547                #calculate coefficient
548                beta1,beta2,delta1,delta2,delta3 = self.coef(self.r)
549                #determine error vector
550                err_vec = self.a * (zk[2*n:] - (delta1*self.zkm1[:n] \
551                        + delta2*self.zkm1[n:2*n] + delta3*self.zkm1[2*n:] \
552                        + beta1 * zk[:n] + beta2 * zk[n:2*n]))
553                #normalize error
554                err = self._norm(err_vec,sc)
555                return err
556
557        def _adjust_h(self,err_mag,TOL = 1e-5,step_accept = True,itt = 0):
558            #function to adjust the step size
559            fac = 0.9 * (2*self.maxit+1)/(2*self.maxit+itt)
560            return fac * self.h * err_mag ** (-1/4)
561
562
563 if __name__ == '__main__':
564     #-----example usage-----
565     from assimulo.ode import Explicit_Problem
566     import matplotlib.pyplot as plt
567
568     #define the problem
569     g : float = 9.81
570     l : float = 0.5
571     theta : float = np.pi/2.
572     def rhs(t,y):
573         return np.array([
574             y[1],
575             np.sin(y[0])*g/l
576             ])
577
578     def J(t,y):
579         return np.array([
580             [0, 1],
581             [np.cos(y[0])*g/l,0]
582             ])
583
584     #--helper functions for calculating period
585     def arith_geo_mean(a,b):
586         while True:
587             a,b = (a+b)/2, np.sqrt(a*b)
588             yield a,b
589
590     from itertools import islice
591     def elliptic_integral(k,tol=1e-5,maxiter=100):
592         a_0,b_0 = 1.,np.sqrt(1-k**2)
593         for a,b in islice(arith_geo_mean(a_0,b_0),maxiter):
594             if abs(a-b) < tol:
595                 return np.pi /(2*a)
596         else:
597             raise Exception('Algorithm did not converge')
```

37

```
598
599     periode = 4*np.sqrt(l/g) * elliptic_integral(np.sin(theta/2),tol=1e-10)
600
601     #----Initialize Problem
602     pend_model = Explicit_Problem(rhs,y0=np.array([theta,0.]))
603     pend_model.name = 'Pendulum simulation'
604     pend_model.jac = J
605     pend_model.h = 0.01
606
607     #----Initialize Solver
608 #   sim = Radau5ODE(pend_model)
609     sim = Lobatto4ODE(pend_model)
610     sim.h = 0.016
611     sim._do_comp_err = True
612     sim.maxit=10
613     #----Simulate the solver
614     t,y = sim.simulate(periode)
615     #----Show results
616     sim.plot()
617     plt.show()
```

## 7.2 Experiment implementation

The following code is written by Robert Klöfkorn with the help of [11]

```python
1  import numpy
2  import matplotlib
3  matplotlib.rc( 'image', cmap='jet' )
4  from matplotlib import pyplot
5
6  ########################################################
7  ## Assimulo imports
8  ########################################################
9  #import assimulo.solvers as aso
10 from Lobatto_IIIC_Assimulo import Lobatto4ODE
11 import assimulo.ode as aode
12 import assimulo.solvers as aso
13
14 ########################################################
15 ## DUNE imports
16 ########################################################
17 #from dune.grid import structuredGrid as leafGridView
18 from dune.grid import cartesianDomain
19 from dune.alugrid import aluCubeGrid as leafGridView
20 #from dune.alugrid import aluSimplexGrid as leafGridView
21 from dune.common import FieldVector
22 from dune.grid import reader
23 from dune.fem import parameter
24 #from dune.fem.space import dgonb as dgSpace
25 from dune.fem.space import dglagrangelobatto as dgSpace
26 from dune.fem.operator import molGalerkin as molGalerkin
27 from dune.fem.function import uflFunction, integrate
28 from dune.ufl import Constant
29 from ufl import TestFunction, TrialFunction, SpatialCoordinate, triangle, FacetNormal
30 from ufl import dx, ds, grad, div, grad, dot, inner, sqrt, exp, conditional, sin, cos
31 from ufl import as_vector, avg, jump, dS, CellVolume, FacetArea, atan, pi
32 from dune.femdg.rk import ssp3, euler
33
34 # useAssimulo = False
35 useAssimulo = True
36
37 # 3rd order 4-stage Runge-Kutta (R.Alexander)
38 Stepper = ssp3(4,explicit=False)
39
40 # Explicit/Implicit Euler
41 #Stepper = euler(explicit=False)
42
43 time = Constant(0., "time")
44 dt   = Constant(0.005,"dt")
45
46
47 class SpatialOperator:
48     def __init__(self,form,space,cfl=0.45):
49         # create method of lines Galerkin operator from PDE form
50         self._op = molGalerkin( form )
51         self.cfl = cfl
```

```
52          self.space = space
53          # discrete functions uTmp and vTmp
54          self.uTmp = space.interpolate([0], name='uTmp')
55          self.vTmp = space.interpolate([0], name='vTmp')
56          self.localTimeStepEstimate = [dt.value/self.cfl]
57
58      # v = L[u]
59      def apply(self, u, v):
60          self._op(u,v)
61          self.localTimeStepEstimate = [dt.value/self.cfl]
62          return
63
64      # jacobian of L[\bar{u}]
65      def jacobian(self, ubar):
66          from dune.fem.operator import linear as linearOperator
67          return linearOperator(self._op, ubar=ubar).as_numpy
68
69      # v = L[u]
70      def __call__(self,u,v):
71          self.apply(u,v)
72          return
73
74      # make current simulation time known to operator
75      def setTime(self, t):
76          time.value = t
77          self._t  = t
78
79      def stepTime(self,t0, dt0):
80          global time
81          if hasattr(self._op.model, "time"):
82              print(f"Setting time to {self._t} + {t0 * dt}")
83              self._op.model.time.value = self._t + t0 * dt
84          else:
85              time.value = self._t + t0 * dt
86          # set time to model time if available
87          # since time is not in the form we don't need this here
88          #if hasattr(self._op.model,"time"):
89          #    print("Model has time")
90          #elif hasattr(self._op.model,"t"):
91          #    print("Model has t")
92
93      def applyLimiter(self, u):
94          pass
95
96      ####################################################
97      # rhs function for Assimulo forwarding to apply
98      ####################################################
99      def rhs(self, t, y):
100         """ Function that calculates the right-hand-side. Depending on
101             the problem and the support of the solver, this function has
102             the following input parameters:
103
104             rhs(t,y)      - Normal ODE
105         """
106         ## set time for PDE operator
107         self.setTime(t)
```

```
108          # copy content of y into uTmp
109          self.uTmp.as_numpy[:] = y[:]
110          # apply spatial discretization operator L
111          self.apply(self.uTmp, self.vTmp)
112          # store result in y
113          yn = y.copy()
114          yn[:] = self.vTmp.as_numpy[:]
115          return yn
116
117      def jac(self, t, y, sw=None):
118          # copy content of y into uTmp
119          self.uTmp.as_numpy[:] = y[:]
120          return self.jacobian( self.uTmp ).toarray()
121
122 parameter.append({"fem.verboserank": 0})
123
124 #########################################################
125 ##
126 ## Spatial discretization
127 ##
128 #########################################################
129
130 domain = cartesianDomain([0, 0], [1, 1], [20, 20])
131 gridView = leafGridView(domain, dimgrid=2 )
132 space    = dgSpace(gridView, order=2, storage="fem")
133
134 u    = TrialFunction(space)
135 v    = TestFunction(space)
136 n    = FacetNormal(space)
137 he   = avg( CellVolume(space) ) / FacetArea(space)
138 hbnd = CellVolume(space) / FacetArea(space)
139 x    = SpatialCoordinate(space)
140
141 center  = as_vector([ 0.5,0.5 ])
142 x0 = x[0] - center[0]
143 x1 = x[1] - center[1]
144
145 ux = -4.0*x1
146 uy =  4.0*x0
147
148 # diffusion factor
149 epsilon = 0.001
150
151 # transport direction and upwind flux
152 b =  as_vector([ux,uy])
153 def u0(x,t):
154     sig2 = 0.004
155     sig2PlusDt4 = sig2+(4.0*eps*t)
156     xq = ( x0*cos(4.0*t) + x1*sin(4.0*t)) + 0.25
157     yq = (-x0*sin(4.0*t) + x1*cos(4.0*t))
158     return (sig2/ (sig2PlusDt4) ) * exp (-( xq*xq + yq*yq ) / sig2PlusDt4 )
159
160 # transport direction and upwind flux
161 #b =  as_vector([1,1])
162
163 #def u0(xp,t):
```

41

```
164 #     res = 1.
165 #     for d in range(len(xp)):
166 #         res *= sin(2.*pi *(xp[d] - t))
167 #     return res
168
169
170 # upwind (same as LLF in this case)
171 hatb = (dot(b, n) + abs(dot(b, n)))/2.0
172
173 # diffusion factor
174 eps = Constant(epsilon,"eps")
175
176 # penalty parameter for DG scheme
177 beta = Constant( 10*space.order**2 if space.order > 0 else 1,"beta")
178
179 # exact solution
180 exact = uflFunction(gridView, name="exact", order=3, ufl=u0(x,time))
181
182 # d_t u + div( F(u) - eps grad u) = 0
183 # integration by parts on spatial terms yields
184 aInternal     = inner(eps*grad(u) -b*u, grad(v)) * dx
185 advSkeleton   = jump(hatb*u)*jump(v)*dS \
186                  +(hatb*u + (dot(b,n)-hatb)*exact)*v*ds
187
188 # B(u,v) = \int_\Omega grad(u)*grad(v)
189 #          - \int_\Omega { grad(u) } * [ v ] + [ u ] * { grad(v) }
190 #          + \int_\Gamma \eta*h^-1 [ u ] * [ v ]
191 #
192 # interior skeleton
193 diffSkeleton  = -eps*inner(jump(u,n),   avg(grad(v)))*dS \
194                  -eps*inner(avg(grad(u)),jump(v,n))*dS \
195                  +eps*beta/he*jump(u)*jump(v)*dS
196 # boundary skeleton
197 diffSkeleton += -eps*(u-exact)*dot(grad(v),n)*ds \
198                  -eps*dot(grad(u),n)*v *ds \
199                  +eps*beta/hbnd*(u-exact)*v*ds
200 #rhs = eps*8.*pi*pi* inner(exact,v) * dx
201
202 # minus since we are solving d_t u = L[u]
203 # which leads in the simplest form to
204 # unew = uold + dt * L[uold]
205 # However, in dune-fem we assume that
206 # d_t u + L[u] = 0
207 # therefore we implement -L[u] since we solve
208 # something like d_t u = -L[u]
209 form = -(aInternal + advSkeleton)
210 if abs(epsilon) > 0:
211     form -= diffSkeleton
212 #    form += rhs
213
214 ## Create right hand side operator
215 op = SpatialOperator(form, space, cfl=0.25)
216
217
218 error0 = 0.
219 eoc = 0
```

```python
220
221  # 3 EOC loops
222  for i in range(2):
223      t = 0
224      op.setTime( t )
225      time.value = t
226
227      # interpolate exact solution onto discrete space
228      uh = space.interpolate( exact, name="solution")
229      #uh.plot()
230
231      # write initial data
232      gridView.writeVTK("adv",number=0,celldata=[uh],pointdata=[uh,exact])
233
234      # T
235      endTime = 0.25
236
237      # time derivative
238      if useAssimulo:
239          # uh.plot()
240          y0 = numpy.zeros( uh.size )
241          y0[:] = uh.as_numpy[:]
242
243          AdvDiff = aode.Explicit_Problem(op.rhs, y0, 0.)
244          AdvDiff.name = 'Rotating Pulse'
245          AdvDiff.h = dt.value
246          AdvDiff.jac = op.jac
247
248          ###Choose which solver to use
249          #solver = aso.RungeKutta34(AdvDiff)
250          solver = Lobatto4ODE(AdvDiff)
251          #solver = aso.Radau5ODE(AdvDiff)
252          solver.report_continuously = True
253
254          solver.atol = 1e-5
255          solver.rtol = 1e-5
256          solver.h = 0.0001
257          #Option for our Lobatto method
258          #solver._dynamic_step = False
259
260
261          t, y = solver.simulate(endTime)
262
263
264          # copy last solution back to uh
265          uh.as_numpy[:] = y[-1][:]
266      else:
267          stepper = Stepper(op, cfl=op.cfl)
268
269          uh_n = uh.copy()
270          while t < endTime:
271              stepper( uh, dt.value )
272              t += dt.value
273              op.setTime( t )
274              print(f"time = {t}, dt = {dt.value}")
275
```

```
276    #safe result
277    gridView.writeVTK("adv",number=1,celldata=[uh],pointdata=[uh,exact])
278
279    # compute L2 error
280    error1= numpy.sqrt( integrate(gridView,dot(uh-exact,uh-exact),order=6))
281
282    # plot solution
283    #uh.plot()
284    # exact.plot()
285
286    # compute EOC
287    if i > 0:
288        eoc = [ numpy.log(error1/error0) / numpy.log(0.5) ]
289    print(f"Step {i}: L2-error {error1} | EOC {eoc}")
290    error0 = error1
291
292    # refine grid to half grid size (delta x)
293    gridView.hierarchicalGrid.globalRefine(1)
294    # adjust time step size
295    dt.value *= 0.25
```