

# SPLITTING SCHEMES FOR NONLINEAR PARABOLIC PROBLEMS

TONY STILLFJORD



LUND UNIVERSITY

Faculty of Science  
Centre for Mathematical Sciences  
Numerical Analysis

Numerical Analysis  
Centre for Mathematical Sciences  
Lund University  
Box 118  
SE-221 00 Lund  
Sweden  
<http://www.maths.lu.se/>

Doctoral Theses in Mathematical Sciences 2015:4  
ISSN 1404-0034

ISBN 978-91-7623-252-1 (print)  
ISBN 978-91-7623-253-8 (electronic)  
LUNFNA-1006-2015

© Tony Stillfjord, 2015

Printed in Sweden by Media-Tryck, Lund 2015

# Abstract

This thesis is based on five papers, which all analyse different aspects of splitting schemes when applied to nonlinear parabolic problems. These numerical methods are frequently used when a problem has a natural decomposition into two or more parts, as the computational cost may then be significantly decreased compared to other methods. There are two prominent themes in the thesis; the first concerns convergence order analysis, while the second focuses on structure preservation.

To motivate the first theme, we note that even if a method has been shown to converge it might be that the speed of convergence is arbitrarily slow. As such a method is unusable in practice we see that it is essential to prove convergence orders. However, those studies that present such error analyses in the fully nonlinear setting typically assume more regularity of the solution than what should be expected. In this context, we present a convergence order analysis for a class of splitting schemes which, importantly, does not require any artificial regularity assumptions. This analysis is carried out in the setting of  $m$ -dissipative operators, which includes a large number of interesting problem classes. As demonstrated by the first three papers, the theory can be applied to such diverse problems as nonlinear reaction-diffusion systems, nonlinear parabolic problems with delay, as well as differential Riccati equations.

Within the second theme of structure preservation, an in-depth study of operator-valued differential Riccati equations has been carried out. In such equations it is desirable for a numerical method to produce positive semi-definite approximations. Further, it is essential that an implementation can utilize the problem-inherent property of low rank. As shown in the last three papers, both these features are readily satisfied for various splitting schemes. Since these are additionally less costly than existing comparable methods, they constitute a particularly competitive choice for such problems.



# Populärvetenskaplig sammanfattning

Genom observationer och experiment kan man konstruera modeller för att beskriva otaliga fenomen och aspekter av det universum vi lever i. Den huvudsakliga ingrediensen i en dylik modell är ofta en partiell differentialekvation. Sådana ekvationer kan beskriva så vitt skilda fenomen som t.ex. hur galaxer bildas, hur luftflöden transporteras i atmosfären, hur kemikalier reagerar med varandra, eller hur atomer interagerar på kvantnivå. Ofta används linjära modeller, då de är relativt enkla och har studerats intensivt. I den här avhandlingen intresserar vi oss dock för *ickelinjära* ekvationer. Eftersom många naturliga fenomen är icke linjära kan dessa beskriva verkligheten bättre än linjära ekvationer.

Att *beskriva* ett system med en ekvation är en sak, men för att använda modellen till att förutsäga vad som kommer att hända i olika situationer måste den också *lösas*. De ekvationer som beskriver komplicerade processer likt de ovan nämnda har dock sällan några lösningar som man kan beräkna genom ett ändligt antal matematiska operationer. Istället måste man hitta tillräckligt bra approximationer, uppskattningar, till dessa lösningar. En stor del av numerisk analys handlar om att konstruera, analysera och implementera metoder för att beräkna sådana approximationer. I den här avhandlingen har fokus varit på att analysera, och till viss del implementera, en viss typ av numeriska metoder som kallas *splitting-metoder*.

Idén bakom en splitting-metod är väldigt enkel; dela upp problemet i två eller fler delar och approximerar deras lösningar separat. Använd sen dessa delapproximationer för att konstruera en approximation till lösningen av hela problemet. Om delproblemen är enklare att hantera än det ursprungliga problemet (t.ex. om man exakt vet delproblemets lösningar) så kan detta leda till en kraftig minskning av den datorkraft som krävs.

Att en numerisk metod är snabb betyder dock inte nödvändigtvis att den är bra, utan man måste fråga sig hur noggranna approximationer den producerar. En central frågeställning är hur approximationsfelet beror av hur mycket arbete man investerar. Om mer datorkraft inte resulterar i en bättre approximation så är metoden inte särskilt bra. Huvudtemat i den här avhandlingen har därför varit att visa så kallade konvergensord-

ningar för splittingmetoder. Detta betyder att man t.ex. kan garantera att en dubblering av arbetsinsatsen resulterar i en halvering av felets storlek.

Sådana felanalyser för splittingmetoder har gjorts tidigare i litteraturen, men under antaganden på problemen som i det icke linjära fallet inte är fullt realistiska eller utesluter intressanta fall. Via det nya tillvägagångssättet som beskrivs i den här avhandlingen kan man dock utföra rigorösa felanalyser även under minimala antaganden. Teorin som presenteras är också applicerbar på många olika klasser av icke linjära problem.

Utöver detta huvudspår så har en del av avhandlingen fokuserat på strukturbevarande numeriska metoder. I t.ex. en kemisk reaktion så kan man självklart inte ha negativa koncentrationer av något ämne. En metod borde därför producera approximationer där alla koncentrationer är positiva. En sådan metod bevarar då strukturen positivitet. I den här avhandlingen har så kallade differentiella Riccati-ekvationer studerats, vilkas lösningar uppvisar vissa strukturer som bör bevaras. Splitting-metoder har tidigare inte tillämpats på denna typ av ekvationer, men våra resultat indikerar att de är mycket väl lämpade för att bevara dessa strukturer. Då de även är effektivare än existerande jämförbara metoder så är de inom detta område mycket lovande metoder.

# List of Papers

This thesis is based on the following papers, listed below in chronological order.

- I. E. Hansen, T. Stillfjord, **Convergence of the implicit-explicit Euler scheme applied to perturbed dissipative evolution equations**, *Mathematics of Computation*, 82 (2013), pp. 1975–1985.
- II. E. Hansen, T. Stillfjord, **Implicit Euler and Lie splitting discretizations of nonlinear parabolic equations with delay**, *BIT*, 54 (2014), pp. 673–689.
- III. E. Hansen, T. Stillfjord, **Convergence analysis for splitting of the abstract differential Riccati equation**, *SIAM Journal of Numerical Analysis*, 52 (2014), pp. 3128–3139
- IV. T. Stillfjord, **Low-rank second-order splitting of large-scale differential Riccati equations**,  
To appear in *IEEE Transactions on Automatic Control*, 2015.
- V. T. Stillfjord, **Convergence analysis for the exponential Lie splitting scheme applied to the abstract Riccati equation**,  
Lund University preprints, 2015.

## **Author's contribution**

My contribution to the papers is listed below.

- I. I contributed to the convergence analysis and the design of the experiments. The final implementation is due to me.
- II. I performed the initial study of the delay setting and contributed to the convergence analysis. The implementation and numerical experiments are due to me.
- III. I performed the initial study of the Hilbert–Schmidt framework and did the implementation and numerical experiments. I contributed to the proof of the convergence theorem and the derivation of the low-rank results.
- IV. All the work on this paper is due to me.
- V. All the work on this paper is due to me.



# Acknowledgments

This thesis was not written in a vacuum, and I would like to thank those who have accompanied me along the way.

First and foremost, I would like to thank my supervisor Esquil Hansen, without whom this work would not exist. I am very grateful that you have put up with me and my frequent stupid moments for more than four years. When I grow up one day, I want to be as methodical and meticulous as you.

Next, I want to thank the whole Numerical Analysis group in Lund for providing a very pleasant work environment. In particular, I would like to acknowledge Claus Führer whose friendly personality originally converted me to the field of numerical analysis, and whose knowledge kept me here. I dedicate the abstract to you — this time, you *are* the target audience. Among the senior staff, I am also grateful to Gustaf Söderlind, who knows everything and can give long but coherent explanations of anything. Thank you for always being so supportive and encouraging, regardless of the context.

Let me also thank my fellow PhD students; Erik, Christian, Dara, Fatemeh and Azahar. The corridor would not be the same without our literal open-door policy and entertaining discussions. In particular, let me thank Erik Henningsson for being a better friend to me than I to him, and for answering all *my* stupid questions. I would also like to thank Dara Maghdid for bringing me to Kurdistan, and, equally importantly, also bringing me back.

Finally, I would like to thank my family; Astrid, Lars-Olof and Tove, for always supporting and believing in me. You are much smarter than you think.

Tony Stillfjord  
Lund, 2015



# Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
1.1	Convergence order analysis . . . . .	2
1.2	Structure preservation . . . . .	3
<b>2</b>	<b>Splitting schemes</b>	<b>5</b>
2.1	Basic ideas . . . . .	5
2.2	Finite-dimensional error analysis . . . . .	6
2.3	Infinite-dimensional error analysis . . . . .	7
<b>3</b>	<b>Framework of <math>m</math>-dissipative operators</b>	<b>11</b>
3.1	$m$ -dissipative operators . . . . .	12
3.2	A theorem by Crandall and Liggett . . . . .	13
3.3	Proof sketch . . . . .	15
3.4	Semi-inner products and Hilbert spaces . . . . .	18
<b>4</b>	<b>A new convergence analysis</b>	<b>23</b>
4.1	Unifying idea . . . . .	24
4.2	Applications . . . . .	27
4.2.1	Locally Lipschitz perturbations . . . . .	27
4.2.2	Delay terms . . . . .	29
4.2.3	Differential Riccati equations . . . . .	31
<b>5</b>	<b>A closer look at Riccati equations</b>	<b>33</b>
5.1	Structure-preserving splitting schemes . . . . .	33
5.2	Alternative error analysis . . . . .	34
<b>6</b>	<b>Conclusions</b>	<b>37</b>
	<b>Bibliography</b>	<b>39</b>
	<b>Papers I–V</b>	



# Chapter 1

## Introduction

A vast number of physical phenomena can be described by partial differential equations (PDEs), from large-scale processes such as galaxy formation or atmospheric flows to small-scale processes such as pattern-formation on animal hides or the quantum-mechanical interaction between particles. Constructing, implementing and analyzing methods for approximating the solutions to such complicated problems is a major branch of numerical analysis.

The kind of PDEs we are interested in in this thesis are nonlinear parabolic problems. These contain one time-derivative and two spatial derivatives. In the linear case, the prototypical parabolic problem is the heat equation  $\dot{u} = \Delta u$ , which e.g. models the diffusion of heat throughout a homogeneous medium. In many cases, this model is too simplistic. For example, the diffusion of heat in a plasma or the flow of gas or water in a particular porous medium should rather be modelled by a nonlinear equation of the form  $\dot{u} = \Delta\alpha(u)$ . A typical function  $\alpha$  could be  $\alpha(u) = |u|^r u$  with  $r > 0$ , which means that the diffusive effect increases with the magnitude of  $u$ .

Usually, of course, a physical process involves more than just diffusion. As a concrete example, let us consider a semilinear reaction-diffusion equation given by

$$\dot{u}_k = \Delta u_k + G_k(u_1, \dots, u_s), \quad k = 1, \dots, s.$$

Here,  $u_k$  denotes the concentrations of  $s$  different chemicals that diffuse separately and interact according to the coupling terms  $G_k$ . The interaction could for example involve the production of one substance from a combination of two others in the presence of a third substance acting as a catalyst. It might be that the diffusivity depends on the concentrations of the reactants, or that the process takes place in a non-homogeneous medium. In this case, the semilinear problem turns into a fully nonlinear problem as in the previous paragraph.

By instead considering  $u_k$  to be the population densities of different animal species, the above equation could also be interpreted in the context of population dynamics. Then

the diffusive terms would describe dispersal and migration throughout a habitat, and the reaction term would describe the interaction between predators and prey as well as the population increase or decrease due to births and deaths. In this context, a nonlinear diffusion term naturally captures the desire to avoid overcrowding; when the population density increases, the rate of diffusion increases.

Such problems are particularly well suited for splitting schemes. These numerical methods approximate the solution to a problem by decomposing it into parts and working with each part separately. In the example above, a splitting scheme would iterate between the subproblems

$$\dot{u}_k = \Delta u_k \quad \text{and} \quad \dot{u}_k = G_k(u_1, \dots, u_s).$$

The benefit here is twofold. Firstly, one may use tailored methods for each subproblem. In this case, the diffusive terms are stiff, which requires an implicit method. On the other hand, the reaction term is frequently non-stiff and an explicit method could be used for this subproblem. Secondly, the system decouples, so that one may parallelize the approximation of the subproblems. The end result is a numerical method which is less costly than applying an implicit method to the full problem. The next chapter gives an overview of different splitting schemes.

The main goal of this thesis is to analyze splitting schemes for fully nonlinear parabolic problems. Within this broad statement, two main themes are in focus; convergence order analysis and structure preservation.

## 1.1 Convergence order analysis

Consider the discretization of a PDE by a numerical method. As the examples above indicate, we focus on temporal discretizations and let the spatial part of the problem remain continuous. Thus the temporal results are independent of a subsequent spatial discretization and can serve as a building block for the analysis of full discretizations.

To consider any numerical method at all, clearly it must be convergent. That is, if  $h > 0$  denotes the mesh width of a (equidistant) temporal discretization and  $u^n$  approximates the exact solution at the fixed time  $nh$  then we must have

$$\|u^n - u(nh)\| \rightarrow 0 \quad \text{as } h \rightarrow 0$$

for a suitable norm  $\|\cdot\|$ . However, this convergence could be arbitrarily slow, leading to a method that works in theory but not in practice. It is therefore essential to show convergence with an order, i.e.

$$\|u^n - u(nh)\| \leq Ch^p, \tag{1.1}$$

for positive constants  $C$  and  $p$ .

In this context, one has to carefully consider what abstract framework to use. That is, we ask:

*What vector fields do we allow and in what spaces are they defined?*

One approach is to take a specific problem or problem class, such as the reaction-diffusion equations, and tailor the analysis to a specific space. Another approach, used here, is to identify the crucial properties that different interesting problems have in common, and to employ general arguments to simultaneously show results for all of them. While the former approach might result in e.g. improved error bounds, with the right framework the latter can be surprisingly effective. In this thesis, we consider the powerful setting of  $m$ -dissipative operators, which includes many interesting problem classes. An overview is given in Chapter 3.

As indicated by the literature overview in Chapter 2, most convergence order studies make regularity assumptions on the exact solution. In the fully nonlinear case, these assumptions are frequently excessive, i.e. such regularity cannot be guaranteed. In view of this, the first aim of the thesis can be summarized as:

**Aim.** *To prove convergence orders for splitting schemes applied to fully nonlinear parabolic equations, under no artificial regularity assumptions.*

The main contribution within this area is outlined in Chapter 4. This consists of the material in Papers I-III, which prove convergence orders for several diverse problem classes and demonstrate the benefits of splitting schemes in these contexts. The current exposition unifies the treatments in these papers and also generalizes their results to a certain extent.

## 1.2 Structure preservation

To describe structure preservation, we note that it is often not enough that a method converges, it also has to produce approximations that mimic the features of the exact solution. In the above examples, clearly a chemical concentration or a population density cannot be negative. Thus a numerical approximation that does not preserve positivity is undesirable. Other features one might want to preserve are global invariants such as mass or energy or local features such as area or volume. In this context, we ask:

*What numerical methods are suitable for preserving a given feature of a problem?*

In contrast to the first theme, this is naturally a problem-specific question, and we thus also consider it for a specific problem class, namely the differential Riccati equations. For these problems, the (operator-valued) solutions are positive semi-definite, and it is desirable for a method to preserve this feature. For implementation reasons, it is also essential to preserve the feature of low rank. The second aim of the thesis can thus be summarized as:

**Aim.** *To do an in-depth study of splitting schemes applied to differential Riccati equations with the goal of preserving positivity and low-rank structure.*

The main contribution within this area is outlined in Chapter 5. This contains material from Paper IV and V, which show that exponential splitting schemes are well suited to the given task, in addition to being more efficient than comparable methods.



# Chapter 2

## Splitting schemes

Let us formally consider the equation

$$\dot{u} = (F + G)u, \quad u(0) = u_0, \quad (2.1)$$

where  $F$  is typically a nonlinear diffusion operator and  $G$  is a reaction term. The main idea behind *splitting schemes* is that numerically approximating the solutions to the subproblems

$$\dot{u} = Fu \quad \text{and} \quad \dot{u} = Gu \quad (2.2)$$

can be significantly cheaper and/or easier than for the full problem (2.1). A splitting method will iterate between solving the different subproblems, and combine the approximations to an approximation for the full problem.

In the following, we denote by  $e^{tE}u_0$  the solution to  $\dot{u} = Eu$ ,  $u(0) = u_0$ . This is an extension of the notation commonly used for linear ordinary differential equations (ODEs) and will be motivated in the next chapter. For now, we note that  $e^{tE}$  is a potentially nonlinear operator, and as such does not commute with  $E$ .

### 2.1 Basic ideas

As a first example, consider the exponential Lie splitting. The time-stepping operator for this method is given by

$$S_h = e^{hF}e^{hG},$$

and the approximation  $u^n$  to  $u(nh)$  is given by the recursion formula  $u^{n+1} = S_h u^n$  with  $u^0 = u_0$ . In other words, one step of the method is given by

$$\begin{aligned} \dot{v} &= Gv, & v(0) &= u^n \\ \dot{w} &= Fw, & w(0) &= v(h), \end{aligned}$$

and  $u^{n+1} = w(h)$ . This should be contrasted to the operator  $e^{h(F+G)}$  for the exact solution, which utilizes the full vector field.

By combining the subproblem solutions in better ways, one constructs more accurate methods. For example, the second-order Strang splitting [84] is given by

$$S_h = e^{h/2F} e^{hG} e^{h/2F}.$$

More generally, the time-stepping operators

$$S_h = \prod_{k=1}^m e^{\alpha_k h F} e^{\beta_k h G}$$

where  $\alpha_k$  and  $\beta_k$  are chosen appropriately [38, 52, 69, 74] yield exponential splitting methods of higher order.

The above methods all employ the *exact* solutions to the subproblems. If known analytical expression for these do not exist, one may still use numerical methods to approximate them, as long as their errors are negligible compared to the error introduced by the splitting. Alternatively, one may directly analyze “full” splitting methods such as the (non-exponential) Lie splitting scheme, or the IMEX (implicit explicit) Euler method [52, Chapter IV.4], given by

$$S_h = (I - hF)^{-1}(I - hG)^{-1} \quad \text{and} \quad S_h = (I - hF)^{-1}(I + hG),$$

respectively. These should be contrasted to the implicit Euler and explicit Euler methods given by  $(I - h(F + G))^{-1}$  and  $I + h(F + G)$ . Even if  $G$  is non-stiff, a stiff  $F$  would prevent using explicit Euler on the full problem, while the use of implicit Euler would frequently be costly. In such a case, IMEX-type schemes are competitive since they only employ the implicit method where strictly necessary. Furthermore, as in the reaction-diffusion example in the introduction, the evaluation of the subproblems may often be parallelized.

For an introductory survey of splitting methods, we refer to Hundsdorfer and Verwer [52], which also covers several other kinds of splitting methods. See also the survey article [69]. In view of the first aim of the thesis, to prove convergence orders for splitting schemes, let us now consider what has been done in the literature.

## 2.2 Finite-dimensional error analysis

Consider first the ODE case, i.e. equations given on a finite-dimensional space such as  $\mathbb{R}^n$ . Then order conditions for both one-step and multistep methods such as Runge–Kutta or backward differentiation formula (BDF) methods are presented in e.g. the classic books [39, 40]. These are based on expanding both the exact solution and the numerical approximation in Taylor series and then determining the method coefficients such that

the terms match up to the desired order. A systematic approach for doing such expansions in a reasonably clean fashion, also for high orders, was mainly introduced by Butcher (for Runge-Kutta methods) and is therefore referred to as the B-series approach. See [16, 39] for an overview.

The same approach can be used for exponential splitting methods. To give a concrete elucidating example, consider the exponential Lie splitting scheme in the linear case. Then  $F = A$  and  $G = B$  are matrices in  $\mathbb{R}^{N \times N}$  and we can expand  $e^{hA}$  and  $e^{hB}$  as well as  $e^{h(A+B)}$  in their power series to verify that

$$\begin{aligned} e^{h(A+B)} - e^{hB}e^{hA} &= I + h(A+B) + h^2/2(A^2 + AB + BA + B^2) + \mathcal{O}(h^3) \\ &\quad - (I + hB + h^2/2B^2 + \mathcal{O}(h^3))(I + hA + h^2/2A^2 + \mathcal{O}(h^3)) \\ &= h^2/2(AB - BA) + \mathcal{O}(h^3). \end{aligned}$$

Thus  $e^{h(A+B)}u^n - e^{hB}e^{hA}u^n = \mathcal{O}(h^2)$  as  $A$  and  $B$  are bounded operators, and the scheme is consistent of order  $p = 1$ . Since

$$\|e^{hA}e^{hB}\| \leq e^{h(\|A\| + \|B\|)},$$

the scheme is also stable, and hence it is first-order convergent. In the case that  $F$  and  $G$  are nonlinear, then one may still expand the flows in Taylor series and show a  $\mathcal{O}(h^2)$  local error that now depends on the bounded operators  $\frac{dF}{du}G$  and  $\frac{dG}{du}F$ . See e.g. [52] for further details. For high orders, B-series may again be used, and we refer to [74], see also the related [18].

## 2.3 Infinite-dimensional error analysis

In the parabolic case,  $F$  and  $G$  will typically be unbounded operators on a Banach space. Then the above approach immediately fails, as the Taylor expansions no longer necessarily converge. This can be overcome in different ways. We consider first the case when both  $F$  and  $G$  are linear.

### The linear case

For exponential splitting schemes, it was shown in [45] that the classical orders derived in the ODE setting remain valid under certain regularity assumptions on the solution to the full problem. Essentially, for order  $p$ , one needs to be able to apply a combination of  $p + 1$   $F$ 's and  $G$ 's to  $u(t)$  and have the result be uniformly bounded over the integration interval. That is, for first-order convergence e.g.  $F^2u(t)$  needs to be uniformly bounded for  $t \in [0, t_{\text{end}}]$ . Under similar assumptions, [25] considers partially time-dependent equations and shows convergence orders for IMEX-type multistep methods. For an introductory reading on analysis of other methods for linear parabolic problems,

such as Runge–Kutta or multistep methods, see e.g. [53, 61, 93]. These also consider full discretizations, i.e. both temporal and spatial discretizations.

Sufficient regularity of the solution to the full problem may be guaranteed by assuming high regularity of the initial condition. This is e.g. used in [85, p.57] to prove first-order convergence of the implicit Euler method. Similar assumptions are made in [44], where convergence orders for various first- and second-order splitting schemes are established for a wide class of operators. In [54], regularity assumptions are avoided but instead  $G$  is supposed to be bounded, and terms such as  $FG$  should be bounded by powers of  $F$ . In this setting, the classic convergence orders for exponential Lie and Strang splitting are shown to remain valid. Under similar boundedness assumptions, W-methods [40, IV.7] are shown to converge with orders in [76].

### The semilinear case

The semilinear case, i.e.  $F$  is linear and  $G$  is nonlinear, needs additional requirements on  $G$ . If the solution is sufficiently regular and  $G$  is smooth then the implicit Euler and Crank–Nicolson methods converge with orders [93]. Under the same regularity assumption but with  $G$  only locally Lipschitz, convergence orders for Runge–Kutta methods were established in [67].

A Lipschitz-type assumption on  $G$  was again used to prove classical convergence orders for several different splitting schemes in [2, 3], under the assumption that the solution is smooth. Other splitting methods for reaction-diffusion problems, where  $F = \Delta$ , are treated in e.g. [29] and [34] under the assumption that  $G$  is a scalar function and sufficiently differentiable. Splitting methods for the time-dependent case and also with more than two operators was considered in [86] under the assumption that  $G$  is dissipative, but without showing convergence orders. IMEX-type schemes are considered in [26, 60] along with spatial discretizations, and are shown to converge with orders under smoothness or local Lipschitz assumptions on  $G$ , but without regularity assumptions on the initial condition.

Extensions to the quasilinear case of operators of the form  $Fu = \alpha(u)Au$  or  $Fu = B(u)u$ , with  $A$  and  $B(u)$  linear, have been considered in e.g. [37, 68] for Runge–Kutta methods, in [62] for multistep methods and in [93, Chapter 13] for methods combined with Galerkin spatial discretizations.

### The nonlinear case

The literature in the fully nonlinear case is more sparse, even when considering methods that are not of splitting-type. In the case that a linearization can be done, semilinear techniques can be employed. This approach has been used to prove convergence orders for the implicit Euler scheme [36], Runge–Kutta methods [77], as well as multistep methods [78]. In the variational setting, weak convergence has been shown for multistep

methods in [31, 32] and for Runge–Kutta methods in [33]. Convergence orders for multistep and Runge–Kutta methods applied to  $m$ -dissipative vector fields has been shown in [41, 42], under regularity assumptions on the solution.

In the splitting case, convergence *without* orders was shown for the Lie and exponential Lie splitting schemes in [14] in the setting of dissipative operators on Hilbert spaces. The main tool to extend these results to the Banach space setting is given by [15, 71], which provides stability and consistency criteria for convergence of general time-stepping schemes. These results were e.g. utilized by [64] for the Peaceman- and Douglas–Rachford schemes, and developed further in [46] for several other methods, with applications to quasi-linear dimension splitting. An alternative approach based on viscosity solutions was employed in [55] to demonstrate convergence with an order for the scheme  $e^{hF}(I + hG)$  applied to a class of nonlinear strongly degenerate parabolic equations.

It should also be noted that there has been many studies of splitting methods applied to specific problem classes, e.g. conservation laws [20, 21, 49, 90], convection perturbed by diffusion [56, 57], semilinear Burgers-type problems [50, 51], Schrödinger equations [35, 65, 91, 92] and Navier–Stokes equations [87, 88].



## Chapter 3

# Framework of $m$ -dissipative operators

We consider now the abstract problem

$$\dot{u} = Eu, \quad u(0) = u_0, \quad (3.1)$$

where  $E$  denotes a generic operator that could be  $F$ ,  $G$  or  $F + G$ . We need to specify what kind of operators  $E$  can represent, and in which sense they should be analysed.

For problems dominated by convection, one may utilize the framework of entropy or viscosity solutions, advocated by e.g. Holden et.al. We refer to [17, 22] for an introductory reading on these kinds of solutions with basic existence and uniqueness results. A recent survey of the field, in the context of splitting methods, is given by Holden et.al. [49] and includes an extensive bibliography.

For parabolic problems, a common approach is the variational setting. There, the operators are treated as bounded operators from one space into its dual. For example, the variational treatment of the Dirichlet Laplacian would be  $\Delta : H_0^1(\Omega) \rightarrow H^{-1}(\Omega)$ . We refer to Roubířek [81], Zeidler [96] and Thomée [93]. This approach has natural connections to Galerkin methods, and can easily be adapted to treat time-varying vector fields. However, it is most suited to Hilbert space theory and tends to produce weakly convergent approximations since one of the main techniques is based on compactness.

An alternative approach for parabolic problems is that of  $m$ -dissipative operators, see e.g. Barbu [6, 7] for the nonlinear case or Pazy [79] for the linear case. In contrast to the variational setting, here the operators are considered to be unbounded and defined on a subset of a Banach space. For example, the Dirichlet Laplacian is seen as an operator  $\Delta : H^2(\Omega) \cap H_0^1(\Omega) \subset L^2(\Omega) \rightarrow L^2(\Omega)$ . As we shall see, a main benefit of this approach is the possibility of deriving convergence order results under only minimal regu-

larity assumptions on the initial condition. Additionally, it allows for a natural treatment of nonlinear operators on e.g.  $L^r(\Omega)$  with  $r \neq 2$ .

### 3.1 $m$ -dissipative operators

Let us therefore consider Equation 3.1 as given on the Banach space  $(X, \|\cdot\|)$  with an unbounded and possibly nonlinear operator  $E$ .

**Definition 1.** An operator  $E : \mathcal{D}(E) \subset X \rightarrow X$  is *dissipative* if for all  $h > 0$  and  $u, v \in \mathcal{D}(E)$  it holds that

$$\|(I - hE)u - (I - hE)v\| \geq \|u - v\|.$$

If there is a constant  $M[E] \geq 0$  such that  $E - M[E]I$  is dissipative, we say that  $E$  is *shift-dissipative* with the shift  $M[E]$ . Finally,  $E$  is (shift-)  *$m$ -dissipative*<sup>1</sup> if in addition

$$\mathcal{R}(I - hE) = X$$

for all  $h > 0$  such that  $hM[E] < 1$ .

**Example 1** (Linear  $m$ -dissipative operator). Let  $\Omega$  be a bounded domain in  $\mathbb{R}^d$  with a sufficiently regular boundary. Then the Laplacian,

$$Eu = \Delta u,$$

with homogeneous Dirichlet, Neumann or periodic boundary conditions, is  $m$ -dissipative on  $L^2(\Omega)$ . The dissipativity follows easily from integration by parts (using Proposition 1 on page 19), while the range condition can be shown using e.g. the Lax-Milgram theorem and regularity arguments, see e.g. [13, Chapters 8.4, 9.5] or [79, Chapter 7.2].

**Example 2** (Nonlinear  $m$ -dissipative operators). Again let  $\Omega$  be a bounded domain in  $\mathbb{R}^d$  with a sufficiently regular boundary. Then the nonlinear  $r$ -Laplacian (usually called  $p$ -Laplacian) is given by

$$Eu = \nabla \cdot (|\nabla u|^{r-2} \nabla u)$$

for  $r > 0$ . This operator, along with suitable boundary conditions, is  $m$ -dissipative on  $L^2(\Omega)$ .

---

<sup>1</sup>In the context of Hilbert spaces, the concept of  $m$ -dissipative operators coincides with that of maximal dissipative operators. In the latter setting, the operators are identified with their graphs and treated as sets in  $X \times X$ . The “maximal” then means that the graph cannot be extended to a larger dissipative set. In the Banach space setting, both concepts are still valid, but they are no longer necessarily the same. Hence the new name, where the  $m$  is related to “maximal” but should rather be read as  $I - hE$  having “maximal range”. We will stick with the  $m$ -dissipative formalism throughout.



A second nonlinear example is the porous medium operator, given by

$$Eu = \Delta(|u|^r u)$$

with  $r > 0$ . This operator, with homogeneous Dirichlet boundary conditions, is  $m$ -dissipative on  $L^1(\Omega)$  or  $H^{-1}(\Omega)$ . For both of these operators, the verification of  $m$ -dissipativity is highly nontrivial. We refer to e.g. [7, p.68ff, p.117ff] and [81, p.101ff, p.105ff].

As these examples show, many interesting applications fit into the framework. Let us therefore turn to the consequences of these properties.

### 3.2 A theorem by Crandall and Liggett

If  $E$  is  $m$ -dissipative, then the resolvent

$$R_h := (I - hE)^{-1} : X \rightarrow \mathcal{D}(E) \subset X,$$

where  $h > 0$ , is well defined. By the dissipativity, we further get that

$$\|R_h u - R_h v\| \leq \|u - v\|,$$

i.e. the resolvent is nonexpansive. Together, the range condition and dissipativity thus guarantees that for all  $v \in X$  the equation

$$(I - hE)u = v$$

has a unique solution  $u \in \mathcal{D}(E)$ . From a numerical analysis point of view, this means that the implicit Euler scheme given by

$$u^{n+1} = R_h u^n, \quad u^0 = u_0,$$

is well defined for all initial conditions  $u_0 \in X$  and step sizes  $h$ . With the particular step size  $h = t/n$ , the value  $u^n = R_{t/n}^n u_0$  approximates the solution to Equation (3.1) at time  $t$ . While these approximations are well defined for all  $n$ , it does not follow directly that they converge as  $n$  tends to infinity. However, we will see in Theorem 1 that they do converge when  $u_0 \in \overline{\mathcal{D}(E)}$ .

Let us first note that the  $m$ -part of the  $m$ -dissipativity can be relaxed. If  $\mathcal{R}(I - hE)$  contains  $\mathcal{D}(E)$ , the resolvent is defined and nonexpansive as a mapping from  $\mathcal{D}(E)$  to  $\mathcal{D}(E)$ . This would make the terms  $R_{t/n}^n u_0$  well defined for  $u_0 \in \mathcal{D}(E)$ . However, even if the sequence converges as we let  $n$  tend to infinity, we can only guarantee that the limit ends up in  $\overline{\mathcal{D}(E)}$ . This motivates the following range condition, which yields a resolvent  $R_h$  mapping  $\overline{\mathcal{D}(E)}$  into itself for all  $h > 0$ :

**Assumption 1.** For all  $h > 0$ , the operator  $E : \mathcal{D}(E) \subset X \rightarrow X$  satisfies

$$\overline{\mathcal{D}(E)} \subset \mathcal{R}(I - hE).$$

It is also sufficient with only shift-dissipativity:

**Assumption 2.** The operator  $E : \mathcal{D}(E) \subset X \rightarrow X$  is shift-dissipative.

Under these assumptions, we can now state the following convergence theorem by Crandall and Liggett [23]:

**Theorem 1.** Let  $E$  satisfy Assumptions 1 and 2. Then the limit

$$e^{tE}u_0 := \lim_{n \rightarrow \infty} R_{t/n}^n u_0$$

exists for all  $u_0 \in \overline{\mathcal{D}(E)}$  and  $t \geq 0$ , and the operator  $e^{tE} : \overline{\mathcal{D}(E)} \rightarrow \overline{\mathcal{D}(E)}$  is nonexpansive for all  $t \geq 0$ . Further, for all  $u_0 \in \overline{\mathcal{D}(E)}$  it holds that  $e^{(t+s)E}u_0 = e^{tE}e^{sE}u_0$  and  $t \mapsto e^{tE}u_0$  is Lipschitz continuous<sup>2</sup> on  $t \geq 0$ . Finally, for  $n \geq 2M[E]t$  and  $u_0 \in \mathcal{D}(E)$  we have

$$\|e^{tE}u_0 - R_{t/n}^n u_0\| \leq 2te^{4M[E]t}n^{-1/2}\|Eu_0\|. \quad (3.2)$$

**Remark 1.** In other words, the implicit Euler scheme converges for all  $u_0 \in \overline{\mathcal{D}(E)}$ . Moreover, if  $u_0 \in \mathcal{D}(E)$  then the method converges with an order,  $p = 1/2$ . The latter unexpected and remarkable fact is one of the cornerstones of our convergence analysis.

In the case of e.g.  $X = \mathbb{R}^N$ , the function  $e^{tE}u_0$  defined in Theorem 1 is a (classical) solution to Equation 3.1. In the current infinite-dimensional setting, this is not necessarily true, but we may still consider it as a more general type of solution:

**Definition 2.** The (Lipschitz) continuous function

$$u(t) := e^{tE}u_0 \quad (3.3)$$

defined by Theorem 1 is called a *mild* solution to Equation (3.1).

This should be contrasted to the following concept of solution to Equation 3.1, called *strong* solutions by e.g. Barbu [7] and Pazy [79].

**Definition 3.** A *strong* solution  $u$  to (3.1) belongs to  $C([0, t_{\text{end}}], X)$  for a given  $t_{\text{end}} > 0$  and is differentiable almost everywhere with  $\dot{u} \in L^1(0, t_{\text{end}}; X)$ . Further, it satisfies  $u(0) = u_0$  and  $\dot{u} = Eu$  almost everywhere.

The concept of a mild solution is further motivated by Theorem II of [23] which asserts that if a mild solution is sufficiently regular, then it is in fact a strong solution:

<sup>2</sup>Thus  $e^{tE}$  is a semigroup of contractions on  $\overline{\mathcal{D}(E)}$ .

**Theorem 2.** *In addition to the assumptions of Theorem 1, let  $E$  be closed and  $u_0 \in \mathcal{D}(E)$ . Then if  $e^{tE}u_0$  is differentiable almost everywhere it is a strong solution to Equation (3.1).*

The following corollary is immediate from the fact that all Lipschitz continuous functions on a reflexive Banach space are differentiable almost everywhere (more generally, on a Banach space with the Radon–Nikodym property):

**Corollary 3.** *Let  $X$  be reflexive. Under the assumptions of Theorem 2,  $e^{tE}u_0$  is a strong solution to Equation (3.1).*

Conversely, if  $E$  is  $m$ -dissipative then every strong solution  $u$  to (3.1) satisfies  $u(t) = e^{tE}u_0$  [7, p.130]. The mild solutions are therefore a natural generalization of the strong solutions.

### 3.3 Proof sketch

We now return to Theorem 1 and its proof, which will be used in the forthcoming convergence analysis of the splitting methods. For the convenience of the reader, we therefore reproduce select parts of it here. We have already seen that the resolvent  $R_h = (I - hE)^{-1}$  is nonexpansive if  $E$  is dissipative. If  $E$  is only shift-dissipative then the resolvent is no longer necessarily nonexpansive. However, as the next Lemma demonstrates, it is still Lipschitz continuous.

**Lemma 4.** *Let  $E$  satisfy Assumption 1 and 2. Then for all  $u, v \in \mathcal{D}(E)$ , positive integers  $n$  and positive  $h$  such that  $hM[E] < 1$ , the resolvent  $R_h$  satisfies*

1.  $\|R_h u - R_h v\| \leq \frac{1}{1-hM[E]} \|u - v\|$ ,
2.  $\|R_h^n u - R_h^{n-1} u\| \leq \frac{h}{(1-hM[E])^n} \|Eu\|$  and
3.  $\|R_h^n u - u\| \leq \frac{nh}{(1-hM[E])^n} \|Eu\|$ .

*The first property also holds for  $u, v \in \overline{\mathcal{D}(E)}$ .*

*Proof.* Property 1 is evident from  $E$  being shift-dissipative. Property 2 follows from property 1 and the identity

$$R_h^n u - R_h^{n-1} u = R_h^n u - R_h^n (I - hE)u.$$

For property 3, we have by property 2 that

$$\|R_h^n u - u\| \leq \sum_{k=1}^n \|R_h^k u - R_h^{k-1} u\| \leq \sum_{k=1}^n \frac{h \|Eu\|}{(1-hM[E])^k} \leq \frac{nh}{(1-hM[E])^n} \|Eu\|,$$

since  $0 < 1 - hM[E] < 1$ . □

We also have the nonlinear resolvent identity:

**Lemma 5.** *Let  $E$  satisfy Assumptions 1 and 2. Then for all  $\lambda, \mu > 0$  and  $v \in \mathcal{R}(I - \lambda E)$ ,*

$$\begin{aligned} \frac{\mu}{\lambda}v + \frac{\lambda - \mu}{\lambda}R_\lambda v &\in \mathcal{R}(I - \mu E) \quad \text{and} \\ R_\lambda v &= R_\mu \left( \frac{\mu}{\lambda}v + \frac{\lambda - \mu}{\lambda}R_\lambda v \right) \end{aligned}$$

*Proof:* Let  $u = R_\lambda v$ . Then

$$\frac{\mu}{\lambda}(I - \lambda E)u + \frac{\lambda - \mu}{\lambda}u = (I - \mu E)u,$$

and the lemma follows immediately.  $\square$

**Lemma 6.** *Let  $E$  satisfy Assumption 1 and 2 and let  $u_0 \in \mathcal{D}(E)$ . Then with  $t > 0$  and sufficiently large positive integers  $n \geq m$  we have*

$$\|R_{t/n}^n u_0 - R_{t/m}^m u_0\| \leq 2te^{4M[E]t} \left( \frac{1}{m} - \frac{1}{n} \right)^{1/2} \|Eu_0\|. \quad (3.4)$$

*Proof sketch.* Following [23], define

$$a_{m,n} = \|R_\mu^n u_0 - R_\lambda^m u_0\|.$$

for  $\lambda \geq \mu > 0$  and  $\lambda M[E] < 1/2$ . Then by rewriting  $R_\lambda^m u_0 = R_\lambda R_\lambda^{m-1} u_0$  and using Lemma 5 we get

$$a_{m,n} \leq \frac{\mu}{\lambda} a_{m-1,n-1} + \frac{\lambda - \mu}{\lambda} a_{m,n-1}. \quad (3.5)$$

As indicated in Figure 3.1, we can solve this recursion in terms of  $a_{k,0}$  and  $a_{0,k}$ :

$$a_{m,n} \leq \sum_{j=0}^{m-1} \alpha^j \beta^{n-j} \binom{n}{j} a_{m-j,0} + \sum_{j=m}^n \alpha^m \beta^{j-m} \binom{j-1}{m-1} a_{0,n-j},$$

where  $\alpha = \frac{\mu}{\lambda}$  and  $\beta = \frac{\lambda - \mu}{\lambda}$ . In view of the third property of Lemma 4, in the case  $M[E] = 0$  this reduces to

$$a_{m,n} \leq \sum_{j=0}^{m-1} \alpha^j \beta^{n-j} \binom{n}{j} (m-j)\lambda \|Eu_0\| + \sum_{j=m}^n \alpha^m \beta^{j-m} \binom{j-1}{m-1} (n-j)\mu \|Eu_0\|,$$

which becomes

$$a_{m,n} \leq \left( \sqrt{(n\mu - m\lambda)^2 + n\mu(\lambda - \mu)} + \sqrt{(n\mu - m\lambda)^2 + m\lambda(\lambda - \mu)} \right) \|Eu_0\|$$

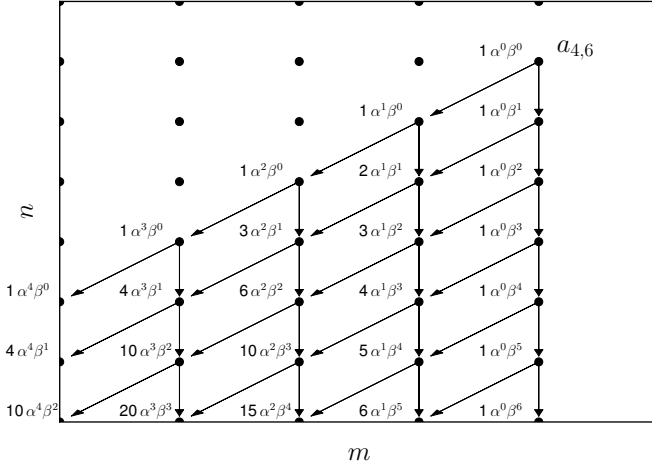


Figure 3.1: Illustrating the Crandall–Liggett proof of Theorem 1. By the recursion formula (3.5), the term  $a_{4,6}$  is first bounded by  $\alpha a_{3,5} + \beta a_{4,5}$ . Each of these terms then give rise to two new terms, indicated by the arrows. Every diagonal movement generates a factor  $\alpha$  and every downward movement generates a factor  $\beta$ . Going down and then diagonally yields a factor  $\alpha\beta$  which is the same result as when going diagonally then down, and we can thus sum these terms. The numbers indicate the total number of the different  $\alpha^i \beta^j a_{k,l}$  terms. We stop when we reach the coordinate axis. In this case, we will have e.g. the terms  $4\alpha^4 \beta a_{0,1}$  and  $15\alpha^2 \beta^4 a_{2,0}$  left.

after using a few combinatorial identities that we omit here, see [23]. The lemma then follows by setting  $\mu = t/n$  and  $\lambda = t/m$ . If  $M[E] > 0$ , we will additionally get terms of the form  $(1 - M[E]\mu)^{-n}$  or  $(1 - M[E]\lambda)^{-n}$ . However, due to the assumption  $M[E]\lambda < 1/2$ , these can be bounded by e.g.  $e^{2M[E]n\mu} = e^{2M[E]t}$ , which yields the factor  $e^{4M[E]t}$  in Equation (3.4).  $\square$

The sequence  $R_{t/n}^n u_0$  is thus a Cauchy sequence, so we can let  $n$  tend to infinity to obtain Equation (3.2). The rest of the proof of Theorem 1 consists of verifying that  $e^{tE}$  actually is a semigroup, and using the Lipschitz continuity to demonstrate that it is defined also on  $\overline{\mathcal{D}(E)}$ . The Lipschitz continuity of  $t \mapsto e^{tE} u_0$  follows from taking  $\mu = t/n$  and  $\lambda = s/n$  and letting  $n$  tend to infinity. We omit these details.

**Remark 2.** One might suspect that the above proof could be modified to show convergence with an order also for the Lie splitting scheme  $S_h = (I - hF)^{-1}(I - hG)^{-1}$  if  $F$  and  $G$  are both  $m$ -dissipative. Then  $S_h$  satisfies all the conclusions of Lemma 4.

However, it does not quite satisfy the resolvent identity. Instead, we get

$$S_\lambda u = S_\mu \left( \frac{\mu}{\lambda} u + \frac{\lambda - \mu}{\lambda} S_\lambda u - \mu(\lambda - \mu) G F S_\lambda u \right),$$

and with  $a_{m,n} = \|S_\mu^n u - S_\lambda^m u\|$  this leads to the modified recursion formula

$$a_{m,n} \leq \frac{\mu}{\lambda} a_{m-1,n-1} + \frac{\lambda - \mu}{\lambda} a_{m,n-1} + \mu(\lambda - \mu) \|G F S_\lambda^m u\|.$$

Solving this recursion yields

$$\begin{aligned} a_{m,n} \leq & \sum_{j=0}^{m-1} \alpha^j \beta^{n-j} \binom{n}{j} a_{m-j,0} + \sum_{j=m}^n \alpha^m \beta^{j-m} \binom{j-1}{m-1} a_{0,n-j} \\ & + \sum_{j=1}^m \frac{1}{\lambda^2} \|G F S_\lambda^j u\| \alpha^{m+1-j} \sum_{k=0}^{j+1} \binom{k+m-j}{k} \beta^{k+1}, \end{aligned}$$

and even assuming that the  $\|G F S_\lambda^j u\|$  terms are uniformly bounded, the  $1/\lambda^2$  factor, essentially  $n^2$ , grows too quickly to achieve a bound similar to those of the other terms. The proof is thus very delicate and specifically tailored for the implicit Euler method.

**Remark 3.** The implicit Euler method is usually said to be a first-order method, i.e. that it is convergent of order  $p = 1$ . The claim in Theorem 1 that we have convergence of order  $p = 1/2$  might therefore seem underwhelming. Indeed, if we take e.g.  $X = \mathbb{R}^N$  or restrict the class of operators further, we do recover first-order convergence. However, in the presented setting and under no additional assumptions, Theorem 1 is *sharp*. This is e.g. demonstrated in [82], where a class of problems is constructed such that the convergence order comes arbitrarily close to  $p = 1/2$ . The result of a numerical verification of this is shown in Figure 3.2, which demonstrates convergence of order  $p = 0.55$ . Convergence orders  $p < 1$  were also observed in Paper I, Example 2, for a splitting method applied to a perturbed version of this problem.

### 3.4 Semi-inner products and Hilbert spaces

Corollary 3 shows that all mild solutions are strong solutions if  $u_0 \in \mathcal{D}(E)$  and  $X$  is reflexive. In general, the more structure  $X$  has, the stronger results can be shown. If the underlying space is a Hilbert space, one may show that  $m$ -dissipative operators are in one-to-one correspondence with the semigroups of contractions. That is, not only does every  $m$ -dissipative operator give rise to a semigroup, but *every* semigroup of contractions arises from an  $m$ -dissipative operator. While this result is interesting from a theoretical viewpoint, from a numerical analysis point of view we always know the operator and

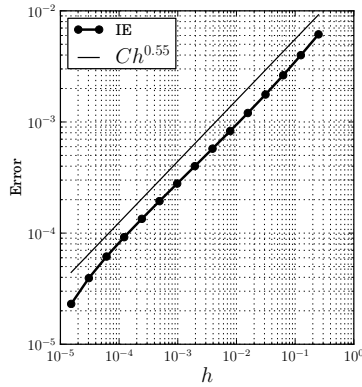


Figure 3.2: The result of a numerical verification of Example 3 in [82]. Here  $X = \ell^2$ ,  $G = 0$  and

$$F(u_1, u_2, \dots, u_{2j-1}, u_{2j}, \dots) = (u_2, -u_1, \dots, ju_{2j}, -ju_{2j-1}, \dots).$$

To provoke low-order convergence behaviour, the initial condition  $u(0) = \{1/j^{1.51}\}_{j=1}^\infty$  is chosen, which belongs to  $\mathcal{D}(F)$  but not to  $\mathcal{D}(F^2)$ . The “spatial discretization” of the problem consists of truncating the series after the first 1000 components, and we compute the error of the implicit Euler method at time  $t = nh = 1$  for different time steps  $h$ . The observed convergence with order  $p = 0.55$  is very close to the theoretical lower bound of 0.5. Compare also Example 2, Paper I.

want to construct or approximate the semigroup. Since the result additionally only holds if we allow *multi-valued* operators, which gives rise to an extra unnecessary level of notational complexity, we do not pursue this line of thought further. See e.g. [70, Theorem 4.20], [24, Theorem A1/A2] or [6, 70] for an extensive literature on similar results.

However, we can employ the extra geometrical features of a Hilbert space to give a useful alternative definition of a dissipative operator.

**Proposition 1.** *Let  $H$  be a Hilbert space with the inner product  $(\cdot, \cdot)$ . Then the operator  $E : \mathcal{D}(E) \subset H \rightarrow H$  is dissipative if and only if*

$$(Eu - Ev, u - v) \leq 0$$

for all  $u, v \in \mathcal{D}(E)$ .

*Proof.* For the “if” part, we observe that

$$\begin{aligned} \|(I - hE)u - (I - hE)v\|^2 &= \|u - v\|^2 + h^2\|Eu - Ev\|^2 \\ &\quad - 2h(Eu - Ev, u - v) \\ &\geq \|u - v\|^2. \end{aligned}$$

To prove the “only if” part, assume that  $\|(I - hE)u - (I - hE)v\| \geq \|u - v\|$  for all  $h > 0$ . Then by the polarization identity

$$\begin{aligned} h(Eu - Ev, u - v) &= \frac{\|(I + hE)u - (I + hE)v\|^2 - \|(I - hE)u - (I - hE)v\|^2}{4} \\ &\leq h/2(Eu - Ev, u - v) + h^2/4\|Eu - Ev\|^2 \end{aligned}$$

so that  $(Eu - Ev, u - v) \leq h/2\|Eu - Ev\|^2$  for all  $h > 0$ . Letting  $h$  tend to zero completes the proof  $\square$

One may define the concept of a dissipative operator in a way similar to the above also in a Banach space  $X$ . Since this is the setting that was used in Paper I we describe it briefly here, but refer to Deimling [28, Chapter 13.1] for a complete exposition. First note that for  $0 \leq \lambda \leq 1$  we have

$$\begin{aligned} \|u + \lambda v + (1 - \lambda)w\| &= \|\lambda(u + v) + (1 - \lambda)(u + w)\| \\ &\leq \lambda\|u + v\| + (1 - \lambda)\|u + w\| \end{aligned}$$

for all  $u, v, w \in X$ . Thus with  $0 < s < t$  we get

$$\|u + sv\| - \|u\| = \left\| u + \frac{s}{t}tv + \left(1 - \frac{s}{t}\right)0 \right\| - \|u\| \leq \frac{s}{t}(\|u + tv\| - \|u\|)$$

so that  $\varphi : t \rightarrow \frac{\|u+tv\| - \|u\|}{t}$  is a monotonically increasing function. Since  $-\|v\| \leq \varphi(t) \leq \|v\|$ , the limits  $\lim_{t \rightarrow 0^+} \varphi(t)$  and  $\lim_{t \rightarrow 0^-} \varphi(t)$  both exist. The following definition therefore makes sense:

**Definition 4.** The *semi-inner products*  $(\cdot, \cdot)_{\pm} : X \times X \rightarrow \mathbb{R}$  are defined by

$$(u, v)_{\pm} = \|v\| \lim_{t \rightarrow 0^{\pm}} \frac{\|v + tu\| - \|v\|}{t}$$

for all  $u, v \in X$ .

One easily confirms that if  $X$  actually is a Hilbert space, then  $(\cdot, \cdot)_{-}$  and  $(\cdot, \cdot)_{+}$  both coincide with its inner product. They also behave much like inner products. For example, we have



- $(u, u)_\pm = \|u\|^2$ ,
- $(\alpha u, \beta v)_\pm = \alpha\beta (u, v)_\pm$  for  $\alpha\beta > 0$ ,
- $(u, w)_\pm + (v, w)_\mp \leq (u + v, w)_\pm \leq (u, w)_\pm + (v, w)_\pm$ , and
- $|(u, v)_\pm| \leq \|u\|\|v\|$ .

See e.g. [28, Proposition 13.1] for these and more properties. We can now give yet another alternative definition of a dissipative operator:

**Proposition 2.** *The operator  $E : \mathcal{D}(E) \subset X \rightarrow X$  is dissipative if and only if*

$$(Eu - Ev, u - v)_- \leq 0$$

for all  $u, v \in \mathcal{D}(E)$ .

*Proof.* If  $(Eu - Ev, u - v)_- \leq 0$  then  $-(Eu - Ev, u - v)_+ \geq 0$ , so

$$\|u - v\| \lim_{t \rightarrow 0^+} \frac{\|u - v - t(Eu - Ev)\| - \|u - v\|}{t} \geq 0.$$

But this means that  $\|u - v - t(Eu - Ev)\| \geq \|u - v\|$  for all sufficiently small  $t > 0$ . Since the function  $t \rightarrow \frac{\|u - tv\| - \|u\|}{t}$  is increasing, the inequality extends to all  $t > 0$ . Conversely, if  $\|u - v - t(Eu - Ev)\| \geq \|u - v\|$  for all  $t > 0$ , then the limit is nonpositive.  $\square$



## Chapter 4

# A new convergence analysis

While the previous chapter provided an answer to the question of what framework to use, in the present chapter we resolve the first main goal of the thesis, namely to prove convergence orders for splitting schemes applied to fully nonlinear parabolic problems. Consider therefore the equation

$$\dot{u} = (F + G)u, \quad u(0) = u_0, \quad (4.1)$$

where both  $F$  and  $G$  are allowed to be nonlinear. As suggested in Chapter 2, the main problem is lack of regularity if one does not make additional assumptions. To illustrate this, consider the equation

$$\dot{u} = \Delta u^3.$$

A class of explicit solutions to this equation is the Barenblatt solutions [9, 94], in one dimension given by

$$u(t, x) = \frac{1}{t^{1/4}} \left( C - \frac{|x|^2}{12t^{1/2}} \right)_+^{1/2},$$

where  $[\cdot]_+ = \max\{\cdot, 0\}$  and  $C$  is a constant. As illustrated in the right plot of Figure 4.1, these solutions have compact support for  $t > 0$ , and the “corners” at the interface are not smooth. We thus do not have any higher-order regularity in space or time. This should be contrasted to the solution of the heat equation with the same initial condition. As is well-known, and illustrated in the left plot of Figure 4.1, the solution immediately (at  $t > 0$ ) becomes infinitely smooth.

At first glance, nothing in this setting suggests that a splitting method (or any numerical method) would converge with any order at all. However, by the remarkable result of Crandall and Liggett in Theorem 1, the implicit Euler scheme converges with order  $p = 1/2$ . It is therefore a reasonable assumption that the same could hold for other methods as well. Demonstrating that this is indeed the case for different splitting methods in

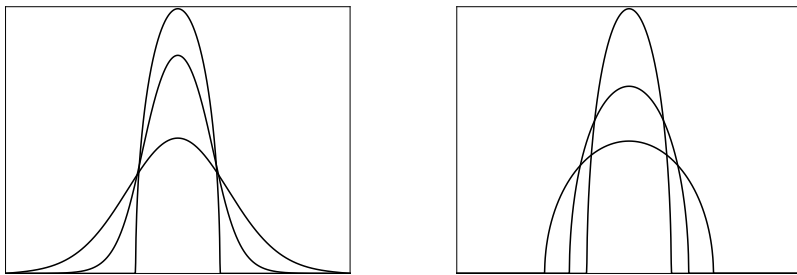


Figure 4.1: Left: The solution to  $\dot{u} = \Delta u$  with a nonsmooth initial condition for three successive times  $t$ . Right: The Barenblatt [9, 94] solution to the PME with the same initial condition at three successive times  $t$ . We see that the linear problem has a smooth solution while the Barenblatt solution is only (Hölder) continuous in both space and time.

several different contexts is the content of Papers I-III. In the Oberwolfach report [48], the slightly different approaches in these papers were unified into one procedure under the assumption that all the operators  $F$ ,  $G$  and  $F + G$  are  $m$ -dissipative. Here, we present the same ideas, but in more detail and generalized to e.g. the alternative range condition from Assumption 1.

## 4.1 Unifying idea

Our aim in this section is to prove convergence orders for the schemes

$$S_h = (I - hF)^{-1}T_{hG},$$

where the method  $T_{hG}$  for the  $G$ -subproblem is yet to be specified. Informally, the idea is that  $F$  is shift-dissipative and satisfies a range condition, and we would like to acquire criteria for convergence that only depend on further properties of  $G$  and  $T_{hG}$ . Depending on the specific problem class, different methods  $T_{hG}$  might be suitable. Formally, we have the following two assumptions. The first one simply guarantees that the method is well defined.

**Assumption 3.** *The domain and range of the operator  $T_{hG}$  satisfy the inclusions*

$$\mathcal{D}(F) \subset \mathcal{D}(T_{hG}) \quad \text{and} \quad \mathcal{R}(T_{hG}) \subset \mathcal{R}(I - hF).$$

The second assumption guarantees the existence of a solution to the full problem. Further, it ensures that both  $(I - hF)^{-1}$  and  $(I - h(F + G))^{-1}$  are Lipschitz continuous and that we can perform the necessary algebraic manipulations.

**Assumption 4.** *The operators  $F$  and  $F + G$  are both shift-dissipative. Further, the domain of  $F + G$  is given by*

$$\mathcal{D}(F + G) = \mathcal{D}(F) \cap \mathcal{D}(G),$$

*and it satisfies the range condition*

$$\overline{\mathcal{D}(F + G)} \subset \mathcal{R}(I - h(F + G)) \quad (4.2)$$

*for all  $h > 0$  such that  $hM[F + G] < 1/2$ .*

Our approach is based on Theorem 1, which guarantees that the implicit Euler scheme, given by the time stepping operator

$$R_h = (I - h(F + G))^{-1},$$

is within a  $\mathcal{O}(h^p)$ -vicinity of the exact solution. Instead of estimating the distance from the splitting approximation to the exact solution, we can thus estimate the distance to the implicit Euler approximation. In the following theorem, and in the rest of the thesis, we denote by  $L[E]$  the Lipschitz constant of an operator  $E : \mathcal{D}(E) \subset X \rightarrow X$ , i.e.

$$L[E] = \sup_{\substack{u, v \in \mathcal{D}(E) \\ u \neq v}} \frac{\|Eu - Ev\|}{\|u - v\|}.$$

**Theorem 7.** *Let Assumptions 3 and 4 be satisfied, let  $u_0 \in \mathcal{D}(F + G)$  be given and suppose that  $h > 0$  satisfies  $h \max(M[F], M[F + G]) < 1/2$ . If  $T_{hG}$  is stable, i.e.,  $L[T_{hG}] \leq 1 + Ch$ , and satisfies the consistency bound*

$$\|(hGR_h + I - T_{hG})R_h^j u_0\| \leq Ch^{1+q}, \quad (4.3)$$

*for all  $j = 0, \dots, n - 1$ , then*

$$\|S_h^n u_0 - u(nh)\| \leq C'(h^p + h^q), \quad 0 \leq nh \leq t_{end},$$

*where  $u$  is the mild solution of (4.1),  $p \in [1/2, 1]$  is the convergence order of the implicit Euler scheme and  $C'$  is a constant which depends on  $t_{end}$  but not on  $n$  or  $h$ .*

*Proof.* By Theorem 1 we have  $\|R_h^n u_0 - u(nh)\| \leq Ch^p$ , so it is enough to estimate  $\|R_h^n u_0 - S_h^n u_0\|$ . We first note that by the assumptions and Lemma 4,  $S_h$  is stable, and we have

$$L[S_h]^{n-j} \leq (1 - hM[F])^{-(n-j)} (1 + Ch)^{n-j} \leq e^{2t_{end}(M[F]+C)}.$$

The theorem thus follows by the following telescopic expansion:

$$\begin{aligned}
\|R_h^n u_0 - S_h^n u_0\| &\leq \sum_{j=1}^n \|S_h^{n-j} R_h^j u_0 - S_h^{n-j+1} R_h^{j-1} u_0\| \\
&\leq \sum_{j=1}^n L[S_h]^{n-j} L[(I - hF)^{-1}] \|((I - hF)R_h - T_{hG})R_h^{j-1} u_0\| \\
&\leq e^{Ct_{\text{end}}} \sum_{j=1}^n \|(hGR_h + I - T_{hG})R_h^{j-1} u_0\|,
\end{aligned}$$

where the  $n$  terms of the sum cancel one power of  $h$  and leaves  $Ch^q$ .  $\square$

**Remark 4.** Under the assumptions of Theorem 7, but with  $u_0 \in \overline{\mathcal{D}(F+G)}$  and  $\overline{\mathcal{D}(F)} \subset \mathcal{D}(T_{hG})$ , the splitting scheme does not necessarily converge with an order. However, as an immediate consequence of Theorem 1 and the above proof, it still converges to the mild solution of (4.1).

**Remark 5.** In addition to the pointwise convergence result of Theorem 7 we may also show convergence with an order in  $L^\infty(0, t_{\text{end}}; X)$  or  $C([0, t_{\text{end}}], X)$ , by interpolating the values  $S_h^n u_0$ . To this end, assume that  $h$  is chosen such that  $t_{\text{end}} = Nh$  with an integer  $N$ , and denote by  $u^h$  and  $v^h$  the piecewise constant and linear interpolants, respectively. These are both functions from  $[0, t_{\text{end}}]$  to  $\mathcal{D}(F+G)$ , given by

$$u^h(t) = S_h^n u_0 \quad \text{and} \quad v^h(t) = \frac{t - nh}{h} S_h^{n+1} u_0 + \frac{(n+1)h - t}{h} S_h^n u_0,$$

if  $t \in [nh, (n+1)h)$ , for  $n = 0, 1, \dots, N-1$ , and  $u^h(t_{\text{end}}) = v^h(t_{\text{end}}) = S_h^N u_0$ .

**Corollary 8.** *Under the same assumptions as Theorem 7, both  $u^h$  and  $v^h$  converge pointwise to the mild solution  $u$  of Equation (4.1). Further,*

$$\|u^h - u\|_{L^\infty(0, t_{\text{end}}; X)} \leq C(h^p + h^q) \quad \text{and} \quad \|v^h - u\|_{C([0, t_{\text{end}}], X)} \leq C(h^p + h^q),$$

where  $p$  is the convergence order of the implicit Euler scheme and  $C$  is a constant dependent on  $t_{\text{end}}$  but not on  $N$ .

The proof is based on the fact that for any  $t \in [nh, (n+1)h)$  we have

$$\|S_h^n u_0 - u(t)\| \leq \|S_h^n u_0 - u(nh)\| + \|u(nh + \tau) - u(nh)\|,$$

where  $\tau \in [0, h)$ . The first term is  $C(h^q + h^p)$  by Theorem 7, while the second is  $Ch$  by the semigroup property of the mild solution, the Lipschitz continuity of  $e^{nh(F+G)}$  and property 3 of Lemma 4 which extends also to  $e^{h(F+G)}$ . Details can be found in Paper I, Corollary 4.4, for the  $u^h$  case, while the proof for  $v^h$  is analogous.

## 4.2 Applications

The stability condition  $L[T_{hG}] \leq 1 + Ch$  and the consistency condition (4.3) have to be verified on a case-by-case basis. This analysis is carried out in Papers I-III. While the stability condition is natural and usually simple to verify, the consistency requires more effort. We summarize the ideas in the following subsections.

### 4.2.1 Locally Lipschitz perturbations

In Paper I, we consider the case of a (locally) Lipschitz operator  $G$ . As  $F$  is typically a diffusion operator, this term will be stiff, while the perturbation  $G$  is often nonstiff. It therefore makes sense to apply the IMEX splitting scheme

$$S_h = (I - hF)^{-1}(I + hG),$$

which only uses the implicit method on the stiff part, and handles the perturbation by the very cheap explicit Euler method. A particular class of problems where this approach yields a high reduction in computational cost is that of systems of the form

$$\dot{u}_k = F_k u_k + G_k(u_1, \dots, u_s),$$

for  $k = 1, \dots, s$ . In this case, applying  $(I - hF)^{-1}$  reduces to applying  $(I - hF_k)^{-1}$  to each component, i.e. the system decouples. The coupling term  $G$  can thus be handled explicitly, and the application of both  $(I - hF)^{-1}$  and  $I + hG$  can be parallelized. This should be contrasted to applying implicit Euler to the full problem, which would require a costly Newton iteration.

To formalize the above paragraph, let  $G : \mathcal{D}(G) \subset X \rightarrow X$  be Lipschitz continuous with Lipschitz constant  $L[G]$ . Further assume that  $F$  is shift- $m$ -dissipative with  $\mathcal{D}(F) \subset \mathcal{D}(G)$ . This directly implies that  $T_{hG} = I + hG$  satisfies Assumption 3. As in Paper I, one can further verify by a fix-point argument that  $F + G$  is shift- $m$ -dissipative on  $\mathcal{D}(F)$  with

$$M[F + G] \leq M[F] + L[G].$$

Assumption 4 is thus also satisfied.

The stability of  $T_{hG}$  follows directly by the Lipschitz continuity of  $G$ , so it only remains to show the consistency. However, for this specific  $T_{hG}$  we get

$$\begin{aligned} \|(hGR_h + I - T_{hG})R_h^j u_0\| &= h\|(GR_h - G)R_h^j u_0\| \\ &\leq hL[G]\|R_h^{j+1} u_0 - R_h^j u_0\| \\ &\leq h^2 L[G](1 - h(M[F] + L[G]))^{-n} \|(F + G)u_0\|, \end{aligned}$$

where the last inequality follows from Lemma 4. As

$$(1 - h(M[F] + L[G]))^{-n} \leq e^{2nh(M[F] + L[G])} \leq e^{2t_{\text{end}}(M[F] + L[G])},$$

we have shown the consistency (4.3) with  $q = 1$ . This means that the IMEX scheme converges with the same order as the implicit Euler scheme.

If  $\mathcal{D}(G) = X$  we can also assume that  $G$  is only *locally* Lipschitz continuous, i.e for all  $r$  with  $0 < r \leq r_0 < \infty$  there are constants  $L_r[G] < \infty$  such that

$$\|Gu - Gv\| \leq L_r[G]\|u - v\|$$

for any  $u, v$  in the ball  $\{w \in X ; \|w\| \leq r\}$ . To analyze such operators, we fix a positive  $r \leq r_0$  and consider the truncation  $G_r : X \rightarrow X$ , given by

$$G_r u = \begin{cases} Gu, & \|u - u_0\| \leq r \\ G\left(\frac{ru}{\|u - u_0\|}\right), & \|u - u_0\| > r \end{cases}.$$

This operator is globally Lipschitz with  $L[G_r] = 2L_r[G]$  (compare also Paper I). Thus for  $0 < h(M[F] + L[G_r]) < 1/2$ , the resolvents

$$R_h := (I - h(F + G))^{-1} \quad \text{and} \quad \tilde{R}_h := (I - h(F + G_r))^{-1}$$

both exist as operators from  $\overline{\mathcal{D}(F + G)} = \overline{\mathcal{D}(F)}$  into itself, and  $\tilde{R}_h$  has a Lipschitz constant satisfying

$$L[\tilde{R}_h] \leq (1 - h(M[F] + L[G_r]))^{-1}.$$

As in Lemma 4, we further get

$$\begin{aligned} \|\tilde{R}_{t/n}^j u_0 - u_0\| &\leq j \frac{t}{n} L[\tilde{R}_{t/n}]^n \|(F + G_r)u_0\| \\ &\leq te^{2t(M[F] + L_r[G])} \|(F + G)u_0\| \leq r \end{aligned}$$

for all  $j = 0, \dots, n$  and for small enough  $t$ , say  $t \leq t_r$ . But this means that

$$\begin{aligned} u_0 &= (I - t/n(F + G_r))\tilde{R}_{t/n}u_0 \\ &= (I - t/n(F + G))\tilde{R}_{t/n}u_0 \\ &= (I - t/n(F + G))(I - t/n(F + G_r))\tilde{R}_{t/n}^2u_0 \\ &= (I - t/n(F + G))^2\tilde{R}_{t/n}^2u_0, \end{aligned}$$

and so on. Continuing the procedure, we finally arrive at

$$u_0 = (I - t/n(F + G))^n \tilde{R}_{t/n}^n u_0.$$

Hence,  $R_{t/n}^n u_0$  exists for all  $n \geq 0$  and

$$R_{t/n}^n u_0 = \tilde{R}_{t/n}^n u_0.$$



The mild solution  $u(t) = \lim_{n \rightarrow \infty} \tilde{R}_{t/n}^n u_0$  to the truncated problem  $\dot{u} = (F + G_r)u$  is therefore also a mild solution to the original problem  $\dot{u} = (F + G)u$  on sufficiently short time intervals. Changing  $r$  to  $r'$  in the construction above only changes the maximal time of existence to  $t_{r'}$  and the mild solutions coincide on  $[0, \min(t_r, t_{r'})]$ . By taking the supremum over  $r$  of  $t_r$ , we thus get a lower bound for the maximal interval of existence  $[0, t_{\text{end}}]$  for the mild solution to  $\dot{u} = (F + G)u$ , which can be stated in terms of the Lambert  $W$  function:

$$t_{\text{end}} \geq \sup_{0 < r \leq r_0} \frac{1}{2(M[F] + L_r[G])} W\left(\frac{2r(M[F] + L_r[G])}{\|(F + G)u_0\|}\right).$$

**Remark 6.** The above presentation of the locally Lipschitz results does not assume that  $X$  is reflexive, as in Paper I. This setting is therefore more general, but means that we can only talk about mild solutions, rather than strong solutions. It should be noted that Example 1 in Paper I is not valid without this extension, as it considers the non-reflexive space  $X = C(\bar{\Omega}) \times C(\bar{\Omega}) \times L^1(\Omega)$ .

## 4.2.2 Delay terms

In Paper II, we consider equations on the form

$$\dot{u}(t) = fu(t) + g\left(u(t-1) + \int_{-1}^0 u(t+\sigma)d\sigma\right),$$

where  $u(t)$  belongs to the Hilbert space  $H$ ,  $g : H \rightarrow H$  is a Lipschitz continuous function, and  $f : \mathcal{D}(f) \subset H \rightarrow H$  is  $m$ -dissipative. Such equations e.g. arise in models of population dynamics that take gestation periods into account, see e.g. [59, 72, 73, 75]. More generally, we consider

$$\dot{u}(t) = fu(t) + g\Phi u_t,$$

where the function  $u_t : \sigma \mapsto u(t + \sigma)$ , is called the *history segment*<sup>1</sup> at  $t$ . For technical reasons, the delay operator  $\Phi$  is given by

$$\Phi\rho = \int_{-1}^0 \rho(\sigma)d\eta(\sigma),$$

where  $\eta$  is either absolutely continuous or of bounded variation with a jump discontinuity at  $-1$ . The typical point delay  $u(t-1)$  corresponds to the operator  $\Phi\rho = \rho(-1)$ , which is realized by  $\eta = \chi_{(-1,0]}$ , where  $\chi$  denotes the characteristic function of an interval. Such equations can be put into a Banach space framework by introducing

$$X = H \times L^r(-1, 0; H; \tau),$$

---

<sup>1</sup>Not to be confused with the time derivative  $\dot{u}$ .

where  $1 \leq r < \infty$  determines the class of initial history segments that can be considered, and where  $\tau$  is a scaling factor, see Paper II. Then with the operators

$$F = \begin{pmatrix} f & 0 \\ 0 & 0 \end{pmatrix} \quad \text{and} \quad G = \begin{pmatrix} 0 & g\Phi \\ 0 & \frac{d}{d\sigma} \end{pmatrix},$$

the original equation turns into

$$\dot{U} = FU + GU,$$

for  $U = (u(t); u_t) \in X$ . In this case, applying the implicit Euler scheme to the full problem requires a costly Newton iteration involving  $f$ . On the other hand, applying the splitting scheme

$$S_h = (I - hF)^{-1}(I - hG)^{-1},$$

turns this iteration into a fixed-point iteration not involving  $f$ . We refer to Paper II for detailed algorithms describing both procedures.

With suitable domains  $\mathcal{D}(F)$  and  $\mathcal{D}(G)$ , and with the definition  $\mathcal{D}(F + G) = \mathcal{D}(F) \cap \mathcal{D}(G)$ , Assumption 4 can be verified by using results from Webb [95]. In fact,  $F$ ,  $G$  and  $F + G$  are all shift- $m$ -dissipative, with  $M[F] = 0$  and with the value of  $M[G] = M[F + G]$  dependent on  $\eta$ . As an additional consequence, the scheme  $T_{hG} = (I - hG)^{-1}$  satisfies both Assumption 3 and the stability condition.

By using the fact that one can find an explicit representation of  $(I - hG)^{-1}$  and a very similar representation of  $(I - h(F + G))^{-1}$ , Paper II shows the consistency (4.3) with  $q = 1 - 1/r$ , albeit only in  $H \times L^1(-1, 0; H; \tau)$  rather than  $H \times L^r(-1, 0; H; \tau)$ . As we are mainly interested in the convergence of the solution itself, rather than of the history segments, this is perfectly fine. In fact, we may see the  $L^1$ -convergence of the history segments as a bonus.

**Remark 7.** By employing semi-inner products rather than inner products, the results in Webb [95], and hence also our results, extend to the case of  $H = X$  being a Banach space rather than only Hilbert. However, if  $X$  is not reflexive we only know that the approximations converge to a mild solution in  $X \times L^r(-1, 0; X; \tau)$ , and we cannot say that the first component of this solution is a mild solution to Equation (4.2.2) as the concept is not applicable.

**Remark 8.** The operator  $f$  is assumed to be densely defined in [95], but this requirement is superfluous. Since the delay operator  $\Phi$  maps history segments with values in  $\overline{\mathcal{D}(f)}$  into  $\mathcal{D}(f)$ , it seems plausible that also the condition  $\mathcal{R}(I - hf) = H$  could be relaxed to  $\mathcal{R}(I - hf) \supset \overline{\mathcal{D}(f)}$ . Indeed, restricting  $\mathcal{D}(G)$  and  $\mathcal{D}(F + G)$  to only contain history segments with values in  $\overline{\mathcal{D}(f)}$  means that Assumption 4 is satisfied. However, under this modification it is no longer clear that  $\mathcal{D}(F) \subset \mathcal{D}(T_{hG})$ , i.e. Assumption 3 might not hold. The issue could be avoided by considering only those delay operators

$\Phi$  that map  $L^r(-1, 0; H; \tau)$  into  $\overline{\mathcal{D}(f)}$ , but unless  $\overline{\mathcal{D}(f)} = H$  this is a very strong assumption.

### 4.2.3 Differential Riccati equations

In Paper III, we consider abstract differential Riccati equations. These are operator-valued equations of the form

$$\dot{P} = A^* \circ P + P \circ A + Q - P \circ P, \quad (4.4)$$

and e.g. arise in the optimal control of PDEs [5, 27, 63]. We use  $P$  instead of the usual  $u$  here mainly to follow the notation commonly used in optimal control, where  $u$  represents a given input function, but also to emphasize that  $P(t)$  is an operator. A typical application is the linear quadratic regulator problem. There, the aim is to minimize the cost functional

$$J(u) = \int_0^{t_{\text{end}}} \|y\|^2 + \|u\|^2 dt,$$

subject to the state and output equations

$$\begin{aligned} \dot{x} &= Ax + u, & x(0) &= x_0, \\ y &= Cx, \end{aligned} \quad (4.5)$$

where  $-A$  is an elliptic differential operator. Given the solution to the DRE (4.4), the optimal input  $u_{\text{opt}}$  is then found in the feedback form  $u_{\text{opt}}(t) = -P(t_{\text{end}} - t)x(t)$ .

We analyze Equation 4.4 in the setting proposed by Temam [89]. This setting is a mix of variational and non-variational techniques in the following sense. It is variational in that we are given a standard Gelfand triple

$$V \hookrightarrow H \cong H^* \hookrightarrow V^*,$$

with real Hilbert spaces  $V$  and  $H$ , and an elliptic, bounded operator  $-A : V \rightarrow V^*$ . It is non-variational in that the operators  $F$  and  $G$  we consider are treated as unbounded operators. To make this precise, we introduce the spaces

$$\mathcal{V} = \mathcal{HS}(H, V) \cap \mathcal{HS}(V^*, H) \quad \text{and} \quad \mathcal{H} = \mathcal{HS}(H, H),$$

where  $\mathcal{HS}(X, Y)$  denotes the Hilbert-Schmidt operators from  $X$  to  $Y$ . Then as observed in [6, 89], we have the new Gelfand triple

$$\mathcal{V} \hookrightarrow \mathcal{H} \cong \mathcal{H}^* \hookrightarrow \mathcal{V}^*,$$

where  $\mathcal{V} = \mathcal{HS}(H, V) \cap \mathcal{HS}(V^*, H)$  and  $\mathcal{V}^*$  is identified with  $\mathcal{HS}(V, H) + \mathcal{HS}(H, V^*)$ . We can now consider Equation (4.4) in  $\mathcal{H}$  by defining

$$\begin{aligned} \mathcal{D}(F) &= \{P \in \mathcal{V}; A^*P + PA - Q \in \mathcal{H}\}, \\ FP &= A^* \circ P + P \circ A - Q \quad \text{for all } P \in \mathcal{D}(F), \end{aligned}$$

and

$$G : \mathcal{H} \rightarrow \mathcal{H}, \quad GP = -P \circ P.$$

Then  $F$  is  $m$ -dissipative if  $Q \in \mathcal{H}$ , and  $G$  is dissipative, see e.g. Barbu [6, Chapter II:3.3] or Temam [89]. However, the range of  $I - hG$  is not all of  $\mathcal{H}$ . This is easily seen from the fact that the scalar equation

$$\alpha + h\alpha^2 = \beta \tag{4.6}$$

has no real solutions unless  $\beta > -1/(4h)$ . This property extends to diagonalizable matrices and thus also to the  $\mathcal{H}$ -subset of finite-rank operators. However, if  $\beta$  in Equation (4.6) is positive, then there is a unique solution for all  $h > 0$ . We therefore consider Equation (4.4) restricted to the closed and convex set

$$\mathcal{C} = \{P \in \mathcal{H} : P = P^* \text{ and } (Pv, v)_H \geq 0 \text{ for all } v \in H\},$$

which is the operator-equivalent of the nonnegative real numbers. This is not a restrictive assumption as the solutions to Equation (4.4) lie in  $\mathcal{C}$  for most interesting applications. One can show that if  $Q \in \mathcal{C}$  then  $F$ ,  $G$  and  $F + G$  with  $\mathcal{D}(F + G) = \mathcal{D}(F) \cap \mathcal{C}$  are all dissipative, and the corresponding resolvents all map  $\mathcal{C}$  into  $\mathcal{C}$  [6, Chapter II:3.3]. Assumption 4 is therefore satisfied for the restrictions  $F|_{\mathcal{C}}$  and  $G|_{\mathcal{C}}$ , i.e. the operators  $F$  and  $G$  restricted to the domains

$$\mathcal{D}(F|_{\mathcal{C}}) = \mathcal{D}(F) \cap \mathcal{C}, \quad \text{and} \quad \mathcal{D}(G|_{\mathcal{C}}) = \mathcal{C}.$$

To verify the other assumptions we need to define  $T_{hG}$ . In this case we choose  $T_{hG} = e^{hG}$ , as the current simple form of  $G$  means that we can solve the corresponding subproblem exactly on  $\mathcal{C}$ . For  $P_0 \in \mathcal{C}$ , we have the explicit representation

$$e^{hG}P_0 = (I + hP_0)^{-1} \circ P_0.$$

This is the main reason for choosing  $F$  to be affine rather than linear. As demonstrated in Paper III, after an appropriate spatial discretization the evaluation of  $e^{hG}P_0$  is essentially free. Additionally, applying implicit Euler to only the affine problem means that no Newton iteration is needed. As a result, the splitting scheme is much less costly than implicit Euler applied to the full problem.

As  $e^{hG}$  maps  $\mathcal{C}$  into  $\mathcal{C}$  by the properties of  $(I - hG)^{-1}$  and since  $\mathcal{C} \subset \mathcal{R}(I - hF|_{\mathcal{C}})$ , Assumption 3 is satisfied. Further,  $L[e^{hG}] \leq 1$  due to the dissipativity of  $G$ , so the scheme is stable. Finally, we can verify the consistency property by observing that  $e^{tG}P_0$  is actually a smooth function of  $t$  for all  $P_0 \in \mathcal{C}$ . We can thus expand it in a Taylor series. This, combined with Lemma 4 and the polynomial form of  $GP$  yields the consistency (4.3) with  $q = 1$ , i.e. the splitting scheme converges with the same order as the implicit Euler scheme. We refer to Paper III for the details.

# Chapter 5

## A closer look at Riccati equations

During the work on Paper III, the importance of differential Riccati equations (DREs) of the form

$$\dot{P} = A^* \circ P + P \circ A + Q - P \circ P$$

became apparent, and the study of these constitute a second central line of work in this thesis. The main difficulty in applying any numerical method to such an equation is the fact that it is operator-valued. This means that a spatial discretization will turn it into a matrix-valued equation. Where a vector-valued equation would have  $N$  unknowns, we thus instead have  $N^2$  unknowns. As the solutions are generally dense matrices, this yields storage problems even for moderate values of  $N$ . It is therefore vital to utilize structural properties of the solution.

### 5.1 Structure-preserving splitting schemes

Taking structural considerations into account has become something of a trend within numerical analysis during the last few decades, in a shift from the focus on “black-box” integrators applicable to “all” problems, to methods tailored for specific problem classes. We refer to the monograph by Hairer, Lubich and Wanner [38] and the recent survey by Christiansen, Munthe-Kaas and Owren [19]. Properties of interest that one would typically want to preserve are e.g. area or volume, first- or second integrals or energy.

In the case of matrix-valued DREs, we would like to utilize the property of low rank. This means that we can factorize the solution  $P \in \mathbb{R}^{N \times N}$  as  $P = ZZ^T$ , where  $Z \in \mathbb{R}^{N \times r}$  with  $r \ll N$ . While there are currently no definite results on when such structure is to be expected, it is frequently observed in practice. Figure 5.1 provides an example of a typical situation. There are partial results in the related context of algebraic Riccati equations [10], and in the context of Lyapunov equations the issue is e.g. discussed in [4, 80, 83].

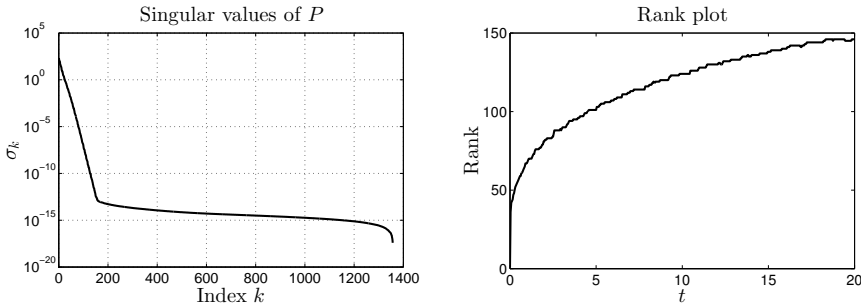


Figure 5.1: Left: The singular values of a typical solution to a matrix-valued DRE. While the problem dimension is  $N = 1357$ , the numerical rank is only about 150. Right: The rank of the splitting approximation of a matrix-valued DRE with larger dimension  $N = 5177$ . Note how the rank stays comparatively very small.

Two further properties we would like to preserve are symmetry and positive semi-definiteness. If  $Q$  and  $P_0$  are symmetric and positive semi-definite, then so is the solution  $P(t)$  for any  $t > 0$  [1]. Clearly, this is automatically fulfilled if a true low-rank implementation is used. However, it was shown in [30] that no linear “one-step” or multistep method of higher order than  $p = 1$  can produce a positive semi-definite approximation. This is exemplified by e.g. the higher-order BDF methods in [11, 12], where the approximation is split into its positive and negative definite parts. Aside from not preserving the positivity, this additionally leads to extra implementation difficulties.

However, as we show in Paper IV, exponential splitting methods allow a true low-rank implementation, and the second-order Strang splitting presents no additional complications as compared to the first-order exponential Lie scheme. This does not contradict the results in [30] because the splitting schemes are *nonlinear* methods.

**Remark 9.** Recently, Koch and Lubich [58] and Lubich and Oseledets [66] have made progress in the related area of dynamical low-rank approximation. In that setting, one fixates a certain rank  $r$  and searches for the best approximation of rank  $r$ . In our case, we search instead for a matrix of lowest rank that yields an error less than a fixed tolerance.

## 5.2 Alternative error analysis

In Paper IV, we also considered more general DREs, of the form

$$\dot{P} = A^* \circ P + P \circ A + Q - P \circ K \circ P, \quad (5.1)$$

where the inclusion of the operator  $K$  yields additional interesting applications in optimal control. Compared to (4.5), we may now consider systems of the form

$$\begin{aligned}\dot{x} &= Ax + Bu, \\ y &= Cx,\end{aligned}$$

which leads to  $K = BB^*$ . However, this modification also means that the nonlinearity  $G$  given by  $GP = -P \circ K \circ P$  is no longer necessarily dissipative, even in the most simple case of  $2 \times 2$ -matrices. As a consequence, the new equation no longer fits into the framework of Paper III, and a different analysis is needed.

Under reasonable assumptions on  $K$ , the operator  $G$  is still locally Lipschitz continuous on  $X = \mathcal{HS}(H, H)$ , and the problem could therefore be put into the framework of Paper I. However, that concerned the splitting scheme  $(I - hF)^{-1}(I + hG)$ . As we know the exact solution  $e^{tG}P_0$  also for this new nonlinearity, using the approximation  $I + hG$  instead seems a waste. Furthermore, as noted in Paper IV, also  $e^{hF}P_0$  can be computed in a straightforward way. Paper V therefore studies the exponential Lie splitting schemes

$$S_h = e^{hF}e^{hG} \quad \text{and} \quad S_h = e^{hG}e^{hF}$$

in the abstract setting of Paper III. We note that in the absence of an operator  $K$ , exponential splitting has also been proposed in [8], but without a convergence order analysis.

The error analysis is now based on comparing the splitting approximation to the exact solution rather than to the implicit Euler approximation. Due to the lack of dissipativity, it is no longer clear that such a solution exists. However, we can temporarily consider Equation (5.1) as semilinear, with the linear part  $L$  given by

$$LP = A^* \circ P + P \circ A,$$

on the domain  $\mathcal{D}(L) = \mathcal{D}(F)$ . If  $G$  is locally Lipschitz then so is the perturbation  $P \mapsto GP + Q$ . The standard results on semilinear equations, e.g. Pazy [79, Section 6.1], therefore yields the existence of a strong solution for sufficiently small time intervals. If  $Q \in \mathcal{D}(L)$  we can additionally show that the perturbation is locally Lipschitz continuous also on the space  $\mathcal{D}(L)$  equipped with the graph norm

$$\|P\|_{\mathcal{D}(L)} = \|P\| + \|LP\|.$$

This means that the strong solution is actually a classical solution [79, Theorem 6.1.7].

Consider now the consistency of the scheme. From the proof of existence, it also follows that the solution  $e^{h(F+G)}P_0$  satisfies the variation-of-constants formula,

$$e^{h(F+G)}P_0 = e^{hL}P_0 + \int_0^h e^{(h-\tau)L} (Ge^{\tau(F+G)}P_0 + Q) d\tau. \quad (5.2)$$

In this formula we can do a first-order Taylor expansion of  $G$ , where the resulting truncation error can be bounded without further regularity assumptions. Similarly, for the

splitting scheme we can do a second-order Taylor expansion of  $e^{hG}$  around  $e^{hF}P_0$ . The difference between the splitting approximation and the exact solution can then be identified as the local error of a first-order quadrature rule, which yields the desired result. This idea originates with Jahnke and Lubich [54] and was used in a similar context to ours in [43], see also the recent paper [47]. Finally, to show convergence of order  $q = 1$ , we must establish also stability of  $e^{hG}$ . However, for sufficiently small time steps this again follows from the local Lipschitz property of  $G$ . Details are provided by Paper V.

**Remark 10.** We note that while this error analysis yields an improved order of convergence compared to the previous results of Paper III, the current approach leads to larger error constants of the form  $e^{tL[G]}$ . This happens since the Lipschitz statements do not take the sign of  $G$  into account;

$$L[-G] = L[G].$$

As a consequence, the theory does not distinguish between solving the equation  $\dot{P} = GP$  in forward or reverse time, and solutions that “blow up” in one temporal direction will be treated as doing so in both directions. By employing dissipativity instead, one can properly treat such cases and verify that a solution exists for all nonnegative times.



# Chapter 6

## Conclusions

In view of the first theme of the thesis, the presented work demonstrates that convergence orders for splitting methods can be shown for fully nonlinear parabolic problems under only minor assumptions on the regularity of the initial condition. While high orders cannot be expected, these results establish a useful baseline and fills the gap in the literature between convergence without orders and convergence with high orders under restrictive assumptions.

We have provided both meta-results such as Theorem 7 as well as verifications of these for explicit problem classes. As demonstrated by the diverse set of applications, the given framework is widely applicable, and further problem classes could be analyzed by applying the meta-results to other perturbations. It should also be reiterated that while the Banach space setting might seem abstract, a main benefit is that the temporal results are independent of subsequent spatial discretizations. The presented results can therefore be used as a cornerstone in the analysis of a full discretization.

Numerical verifications of the convergence results have not been included in this summary, but can be found in the respective papers. These demonstrate the validity of the theory, and also show that the methods are applicable to concrete problems. In addition, they confirm that in many areas splitting schemes constitute a competitive choice.

In view of the second theme of the thesis, the in-depth study of differential Riccati equations demonstrates that splitting schemes are well suited for the preservation of positivity and low rank. This is confirmed by the numerical experiments carried out in Paper III and IV. These are based on algorithms that, while straightforward, constitute the first actual splitting implementations that are applicable to the large-scale setting. As they are less costly than comparable numerical methods, the new convergence order analysis implies that for this problem class, splitting schemes are very promising.



# Bibliography

- [1] H. Abou-Kandil, G. Freiling, V. Ionescu, and G. Jank. *Matrix Riccati equations*. Systems & Control: Foundations & Applications. Birkhäuser, Basel, 2003.
- [2] G. Akrivis and M. Crouzeix. Linearly implicit methods for nonlinear parabolic equations. *Math. Comp.*, 73(246):613–635, 2004.
- [3] G. Akrivis, M. Crouzeix, and C. Makridakis. Implicit-explicit multistep methods for quasilinear parabolic equations. *Numer. Math.*, 82(4):521–541, 1999.
- [4] J. Baker, M. Embree, and J. Sabino. Fast singular value decay for Lyapunov solutions with nonnormal coefficients. arXiv:1410.8741v2 [math.NA], 2015.
- [5] A. V. Balakrishnan. *Applied functional analysis*. Springer, Heidelberg, 1976.
- [6] V. Barbu. *Nonlinear semigroups and differential equations in Banach spaces*. Noordhoff, Leyden, The Netherlands, 1976. Translated from Romanian.
- [7] V. Barbu. *Nonlinear differential equations of monotone types in Banach spaces*. Springer Monographs in Mathematics. Springer, New York, 2010.
- [8] V. Barbu and M. Iannelli. Approximating some nonlinear equations by a fractional step scheme. *Differential Integral Equations*, 6(1):15–26, 1993.
- [9] G. I. Barenblatt. On some unsteady motions of a liquid and gas in a porous medium. *Akad. Nauk SSSR. Prikl. Mat. Meh.*, 16:67–78, 1952.
- [10] P. Benner and Z. Bujanović. On the solution of large-scale algebraic Riccati equations by using low-dimensional invariant subspaces. *MPI Magdeburg Preprints MPIMD/14-15*, 2014.
- [11] P. Benner and H. Mena. BDF methods for large-scale differential Riccati equations. In B. De Moor, B. Motmans, J. Willems, P. Van Dooren, and V. Blondel, editors, *Proc. of Sixteenth International Symposium on: Mathematical Theory of Network and Systems*, 2004.

- [12] P. Benner and H. Mena. Numerical solution of the infinite-dimensional LQR-problem and the associated differential Riccati equations. *MPI Magdeburg Preprint MPIMD/12-13*, 2012.
- [13] H. Brezis. *Functional analysis, Sobolev spaces and partial differential equations*. Springer, New York, 2011.
- [14] H. Brezis and A. Pazy. Semigroups of nonlinear contractions on convex sets. *J. Funct. Anal.*, 6:237–281, 1970.
- [15] H. Brezis and A. Pazy. Convergence and approximation of semigroups of nonlinear operators in Banach spaces. *J. Funct. Anal.*, 9(1):63–74, 1972.
- [16] J. C. Butcher. *The numerical analysis of ordinary differential equations*. Wiley, Chichester, 1987.
- [17] J. Carrillo. Entropy solutions for nonlinear degenerate problems. *Arch. Ration. Mech. Anal.*, 147(4):269–361, 1999.
- [18] F. Casas and A. Murua. An efficient algorithm for computing the Baker–Campbell–Hausdorff series and some of its applications. *J. Math. Phys.*, 50(3):033513, 23, 2009.
- [19] S. H. Christiansen, H. Z. Munthe-Kaas, and B. Owren. Topics in structure-preserving discretization. *Acta Numer.*, 20:1–119, 2011.
- [20] J.-M. Coron. Formules de Trotter pour une équation d’évolution quasilinéaire du premier ordre. *J. Math. Pures Appl. (9)*, 61(1):91–112, 1982.
- [21] M. Crandall and A. Majda. The method of fractional steps for conservation laws. *Numer. Math.*, 34(3):285–314, 1980.
- [22] M. G. Crandall, H. Ishii, and P.-L. Lions. User’s guide to viscosity solutions of second order partial differential equations. *Bull. Amer. Math. Soc. (N.S.)*, 27(1):1–67, 1992.
- [23] M. G. Crandall and T. M. Liggett. Generation of semi-groups of nonlinear transformations on general Banach spaces. *Am. J. Math.*, 93(2):pp. 265–298, 1971.
- [24] M. G. Crandall and A. Pazy. Semi-groups of nonlinear contractions and dissipative sets. *J. Funct. Anal.*, 3:376–418, 1969.
- [25] M. Crouzeix. Une méthode multipas implicite-explicite pour l’approximation des équations d’évolution paraboliques. *Numer. Math.*, 35(3):257–276, 1980.
- [26] M. Crouzeix and V. Thomée. On the discretization in time of semilinear parabolic equations with nonsmooth initial data. *Math. Comp.*, 49(180):359–377, 1987.

- [27] R. F. Curtain and A. J. Pritchard. *Infinite dimensional linear systems theory*, volume 8 of *Lecture Notes in Control and Information Sciences*. Springer, Berlin, 1978.
- [28] K. Deimling. *Nonlinear functional analysis*. Springer, Berlin, 1985.
- [29] S. Descombes and M. Ribot. Convergence of the Peaceman-Rachford approximation for reaction-diffusion systems. *Numer. Math.*, 95(3):503–525, 2003.
- [30] L. Dieci and T. Eirola. Positive definiteness in the numerical solution of Riccati differential equations. *Numer. Math.*, 67(3):303–313, 1994.
- [31] E. Emmrich. Convergence of the variable two-step BDF time discretisation of nonlinear evolution problems governed by a monotone potential operator. *BIT*, 49(2):297–323, 2009.
- [32] E. Emmrich. Two-step BDF time discretisation of nonlinear evolution problems governed by monotone operators with strongly continuous perturbations. *Comput. Methods Appl. Math.*, 9(1):37–62, 2009.
- [33] E. Emmrich and M. Thalhammer. Stiffly accurate Runge–Kutta methods for nonlinear evolution problems governed by a monotone operator. *Math. Comp.*, 79(270):785–806, 2010.
- [34] E. Faou. Analysis of splitting methods for reaction-diffusion problems using stochastic calculus. *Math. Comp.*, 78(267):1467–1483, 2009.
- [35] L. Gauckler and C. Lubich. Splitting integrators for nonlinear Schrödinger equations over long times. *Found. Comput. Math.*, 10(3):275–302, 2010.
- [36] C. González, A. Ostermann, C. Palencia, and M. Thalhammer. Backward Euler discretization of fully nonlinear parabolic problems. *Math. Comp.*, 71(237):125–145 (electronic), 2002.
- [37] C. González and C. Palencia. Stability of Runge–Kutta methods for quasilinear parabolic problems. *Math. Comp.*, 69(230):609–628, 2000.
- [38] E. Hairer, C. Lubich, and G. Wanner. *Geometric numerical integration*, volume 31 of *Springer Series in Computational Mathematics*. Springer, Berlin, second edition, 2006.
- [39] E. Hairer, S. P. Nørsett, and G. Wanner. *Solving ordinary differential equations. I*, volume 8 of *Springer Series in Computational Mathematics*. Springer, Berlin, second edition, 1993.
- [40] E. Hairer and G. Wanner. *Solving ordinary differential equations. II*, volume 14 of *Springer Series in Computational Mathematics*. Springer, Berlin, second edition, 1996.

- [41] E. Hansen. Convergence of multistep time discretizations of nonlinear dissipative evolution equations. *SIAM J. Numer. Anal.*, 44(1):55–65 (electronic), 2006.
- [42] E. Hansen. Runge–Kutta time discretizations of nonlinear dissipative evolution equations. *Math. Comp.*, 75(254):631–640 (electronic), 2006.
- [43] E. Hansen, F. Kramer, and A. Ostermann. A second-order positivity preserving scheme for semilinear parabolic problems. *Appl. Numer. Math.*, 62(10):1428–1435, 2012.
- [44] E. Hansen and A. Ostermann. Dimension splitting for evolution equations. *Numer. Math.*, 108(4):557–570, 2008.
- [45] E. Hansen and A. Ostermann. Exponential splitting for unbounded operators. *Math. Comp.*, 78(267):1485–1496, 2009.
- [46] E. Hansen and A. Ostermann. Dimension splitting for quasilinear parabolic equations. *IMA J. Numer. Anal.*, 30(3):857–869, 2010.
- [47] E. Hansen and A. Ostermann. High-order splitting schemes for semilinear evolution equations. Preprint, 2015.
- [48] E. Hansen and T. Stillfjord. Splitting of dissipative evolution equations. *Oberwolfach Rep.*, 11:814–816, 2014.
- [49] H. Holden, K. H. Karlsen, K.-A. Lie, and N. H. Risebro. *Splitting methods for partial differential equations with rough solutions*. EMS Series of Lectures in Mathematics. EMS, Zürich, 2010.
- [50] H. Holden, K. H. Karlsen, N. H. Risebro, and T. Tao. Operator splitting for the KdV equation. *Math. Comp.*, 80(274):821–846, 2011.
- [51] H. Holden, C. Lubich, and N. H. Risebro. Operator splitting for partial differential equations with Burgers nonlinearity. *Math. Comp.*, 82(281):173–185, 2013.
- [52] W. Hundsdorfer and J. Verwer. *Numerical solution of time-dependent advection-diffusion-reaction equations*, volume 33 of *Springer Series in Computational Mathematics*. Springer, Berlin, 2003.
- [53] A. Iserles. *A first course in the numerical analysis of differential equations*. Cambridge Texts in Applied Mathematics. Cambridge Univ. Press, Cambridge, second edition, 2009.
- [54] T. Jahnke and C. Lubich. Error bounds for exponential operator splittings. *BIT*, 40(4):735–744, 2000.

- [55] E. R. Jakobsen and K. H. Karlsen. Convergence rates for semi-discrete splitting approximations for degenerate parabolic equations with source terms. *BIT*, 45(1):37–67, 2005.
- [56] K. H. Karlsen and N. H. Risebro. An operator splitting method for nonlinear convection-diffusion equations. *Numer. Math.*, 77(3):365–382, 1997.
- [57] K. H. Karlsen and N. H. Risebro. Corrected operator splitting for nonlinear parabolic equations. *SIAM J. Numer. Anal.*, 37(3):980–1003 (electronic), 2000.
- [58] O. Koch and C. Lubich. Dynamical low-rank approximation. *SIAM J. Matrix Anal. Appl.*, 29(2):434–454, 2007.
- [59] Y. Kuang. *Delay differential equations with applications in population dynamics*, volume 191 of *Mathematics in Science and Engineering*. Academic Press Inc., Boston, MA, 1993.
- [60] S. Larsson. Nonsmooth data error estimates with applications to the study of the long-time behavior of finite element solutions of semilinear parabolic problems. Technical Report 36, Dept. of Math., Chalmers University of Technology, 1992.
- [61] S. Larsson and V. Thomée. *Partial differential equations with numerical methods*, volume 45 of *Texts in Applied Mathematics*. Springer, Berlin, 2003.
- [62] M.-N. Le Roux. Méthodes multipas pour des équations paraboliques non linéaires. *Numer. Math.*, 35(2):143–162, 1980.
- [63] J.-L. Lions. *Optimal control of systems governed by partial differential equations*. Springer, New York, 1971.
- [64] P.-L. Lions and B. Mercier. Splitting algorithms for the sum of two nonlinear operators. *SIAM J. Numer. Anal.*, 16(6):964–979, 1979.
- [65] C. Lubich. On splitting methods for Schrödinger–Poisson and cubic nonlinear Schrödinger equations. *Math. Comp.*, 77(264):2141–2153, 2008.
- [66] C. Lubich and I. V. Oseledets. A projector-splitting integrator for dynamical low-rank approximation. *BIT*, 54(1):171–188, 2014.
- [67] C. Lubich and A. Ostermann. Runge–Kutta methods for parabolic equations and convolution quadrature. *Math. Comp.*, 60(201):105–131, 1993.
- [68] C. Lubich and A. Ostermann. Runge–Kutta approximation of quasi-linear parabolic equations. *Math. Comp.*, 64(210):601–627, 1995.
- [69] R. I. McLachlan and G. R. W. Quispel. Splitting methods. *Acta Numer.*, 11:341–434, 2002.

- [70] I. Miyadera. *Nonlinear semigroups*, volume 109 of *Translations of Mathematical Monographs*. American Mathematical Society, Providence, RI, 1992. Translated from Japanese.
- [71] I. Miyadera and S. Ôharu. Approximation of semi-groups of nonlinear operators. *Tôhoku Math. J. (2)*, 22:24–47, 1970.
- [72] J. D. Murray. *Mathematical biology. I*, volume 17 of *Interdisciplinary Applied Mathematics*. Springer, New York, third edition, 2002.
- [73] J. D. Murray. *Mathematical biology. II*, volume 18 of *Interdisciplinary Applied Mathematics*. Springer, New York, third edition, 2003.
- [74] A. Murua and J. M. Sanz-Serna. Order conditions for numerical integrators obtained by composing simpler integrators. *R. Soc. Lond. Philos. Trans. Ser. A Math. Phys. Eng. Sci.*, 357(1754):1079–1100, 1999.
- [75] A. Okubo. *Diffusion and ecological problems: mathematical models*, volume 10 of *Biomathematics*. Springer, Berlin, 1980.
- [76] A. Ostermann. Stability of  $W$ -methods with applications to operator splitting and to geometric theory. *Appl. Numer. Math.*, 42(1-3):353–366, 2002.
- [77] A. Ostermann and M. Thalhammer. Convergence of Runge–Kutta methods for nonlinear parabolic equations. *Appl. Numer. Math.*, 42(1-3):367–380, 2002.
- [78] A. Ostermann, M. Thalhammer, and G. Kirlinger. Stability of linear multistep methods and applications to nonlinear parabolic problems. *Appl. Numer. Math.*, 48(3-4):389–407, 2004.
- [79] A. Pazy. *Semigroups of linear operators and applications to partial differential equations*, volume 44 of *Applied Mathematical Sciences*. Springer, New York, 1983.
- [80] T. Penzl. Eigenvalue decay bounds for solutions of Lyapunov equations: the symmetric case. *Syst. Control Lett.*, 40(2):139–144, 2000.
- [81] T. Roubíček. *Nonlinear partial differential equations with applications*, volume 153 of *International Series of Numerical Mathematics*. Birkhäuser/Springer, Basel, second edition, 2013.
- [82] J. Rulla. Error analysis for implicit approximations to solutions to Cauchy problems. *SIAM J. Numer. Anal.*, 33(1):68–87, 1996.
- [83] D. C. Sorensen and Y. Zhou. Bounds on eigenvalue decay rates and sensitivity of solutions to Lyapunov equations. Technical Report 02-07, Dept. of Comp. Appl. Math., Rice Univ., Houston, TX, 2002.



- [84] G. Strang. On the construction and comparison of difference schemes. *SIAM J. Numer. Anal.*, 5:506–517, 1968.
- [85] H. Tanabe. *Equations of evolution*, volume 6 of *Monographs and Studies in Mathematics*. Pitman, Boston, Mass.-London, 1979.
- [86] R. Temam. Sur la stabilité et la convergence de la méthode des pas fractionnaires. *Ann. Mat. Pura Appl. (4)*, 79:191–379, 1968.
- [87] R. Temam. Sur l’approximation de la solution des équations de Navier-Stokes par la méthode des pas fractionnaires. I. *Arch. Ration. Mech. Anal.*, 32:135–153, 1969.
- [88] R. Temam. Sur l’approximation de la solution des équations de Navier-Stokes par la méthode des pas fractionnaires. II. *Arch. Ration. Mech. Anal.*, 33:377–385, 1969.
- [89] R. Temam. Sur l’équation de Riccati associée à des opérateurs non bornés, en dimension infinie. *J. Funct. Anal.*, 7:85–115, 1971.
- [90] Z. H. Teng. On the accuracy of fractional step methods for conservation laws in two dimensions. *SIAM J. Numer. Anal.*, 31(1):43–63, 1994.
- [91] M. Thalhammer. High-order exponential operator splitting methods for time-dependent Schrödinger equations. *SIAM J. Numer. Anal.*, 46(4):2022–2038, 2008.
- [92] M. Thalhammer. Convergence analysis of high-order time-splitting pseudospectral methods for nonlinear Schrödinger equations. *SIAM J. Numer. Anal.*, 50(6):3231–3258, 2012.
- [93] V. Thomée. *Galerkin finite element methods for parabolic problems*, volume 25 of *Springer Series in Computational Mathematics*. Springer, Berlin, second edition, 2006.
- [94] J. L. Vázquez. *The porous medium equation*. Oxford Math. Monogr. Oxford Univ. Press, Oxford, 2007.
- [95] G. F. Webb. Functional differential equations and nonlinear semigroups in  $L^p$ -spaces. *J. Differential Equations*, 20(1):71–89, 1976.
- [96] E. Zeidler. *Nonlinear functional analysis and its applications. III/B*. Springer, New York, 1990.