# Analytical and numerical approximations to highly oscillatory solutions of nonlinear Friedrichs systems

Tobias Jahnke, Johanna Mödl

KARLSRUHE INSTITUTE OF TECHNOLOGY

CRC 1173

## Participating universities

UNIVERSITÄT BONN

Universität Stuttgart

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

TU
WIEN

## Funded by

DFG

# Analytical and numerical approximations to highly oscillatory solutions of nonlinear Friedrichs systems *

Tobias Jahnke†      Johanna Mödl†

October 28, 2024

### Abstract

We consider semilinear Friedrichs systems which model high-frequency wave propagation in dispersive media. Typical solutions oscillate in time and space with frequency of $\mathcal{O}(\varepsilon^{-1})$ and have to be computed on time intervals of length of $\mathcal{O}(\varepsilon^{-1})$, where $\varepsilon \ll 1$ is a small positive parameter. For such problems, we present an approach which combines analytical approximations with tailor-made time integration. First, we replace the original problem by a fine-tuned modification of the classical slowly varying envelope approximation and prove that the corresponding error is only of $\mathcal{O}(\varepsilon^2)$. The resulting system of partial differential equations has the advantage that solutions do not oscillate in space, but still in time. For this system, we devise a novel time integrator and prove first-order convergence uniformly in $\varepsilon$. Essential to this is the careful analysis of interactions between oscillatory and non-oscillatory parts of the solution, which are identified by suitable projections.

**Keywords:** nonlinear Friedrichs system, high-frequency wave propagation, spatio-temporal oscillations, slowly varying envelope approximation, time integration, error bounds

**MSC classification codes:** 35A35, 35B05, 35B40, 35L45, 35L60, 35Q60, 35Q61, 65M12, 65M15

## 1 Introduction

The construction and analysis of numerical methods for differential equations is particularly difficult if the solution oscillates rapidly in time. With standard schemes, an acceptable accuracy can most often only be achieved if the step size is small compared to the inverse of the highest frequency. This means that the oscillations have to be resolved by a very fine discretization, which causes huge computational costs.

In many cases, however, the efficiency can be substantially improved by exploiting the particular structure of the problem in the construction of the method. In recent years, such tailor-made methods for nonlinear differential equations with highly oscillatory time evolution were proposed and analyzed, e.g., in [1, 2, 6, 9, 10, 11, 12, 13, 16, 18, 19, 20, 22, 23, 24, 27, 26, 30, 33, 34], and many other papers. In this work, we present a new approach for systems of partial differential equations (PDEs) which model high-frequency wave propagation in nonlinear dispersive media. This is particularly challenging, because here the solution oscillates both in time and space, which is not the case in the references cited above.

Our goal is to approximate the vector-valued solution $u : [0, t_{\text{end}}/\varepsilon] \times \mathbb{R}^d \to \mathbb{R}^m$ of the semilinear Friedrichs

system

$$\partial_t u + \mathcal{A}u + \frac{1}{\varepsilon}Eu = \varepsilon T(u,u,u), \qquad t \in (0, t_{\mathrm{end}}/\varepsilon], \, x \in \mathbb{R}^d, \tag{1.1a}$$

$$u(0,x) = \mathrm{e}^{\mathrm{i}\kappa \cdot x/\varepsilon}p_0(x) + \mathrm{e}^{-\mathrm{i}\kappa \cdot x/\varepsilon}\overline{p_0(x)} \tag{1.1b}$$

with a small parameter $\varepsilon > 0$ and differential operator

$$\mathcal{A} = \sum_{j=1}^d A_j \frac{\partial}{\partial x_j}$$

with symmetric matrices $A_1, \ldots, A_d \in \mathbb{R}^{m \times m}$. $E \in \mathbb{R}^{m \times m}$ is skew-symmetric and $T : \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}^m$ trilinear. The right-hand side of the initial condition (1.1b) is a wave packet with a smooth envelope function $p_0 : \mathbb{R}^d \to \mathbb{C}^m$ and a fixed, given wave vector $\kappa \in \mathbb{R}^d \setminus \{0\}$. A well-known example for this type of differential equation is the Maxwell–Lorentz system, which describes the propagation of light in a dispersive Kerr medium, cf. [14, 15, 17, 28, 29] and Section 4.3 below. In order to observe diffractive effects, this PDE system has to be solved on time intervals which are so long that the approximations of geometric optics do not apply.

The problem involves multiple scales in time and space characterized by the parameter $0 < \varepsilon \ll 1$, which appears in the PDE, in the initial data, and in the time interval. Owing to the special form of the initial data and the PDE, the solution is typically a wave packet with a carrier wave, oscillating with frequency proportional to $\varepsilon^{-1}$, and modulated by a smooth, non-oscillatory envelope which moves with group velocity of $\mathcal{O}(1)$. Because of the oscillatory behavior in time and space, the efficient and accurate approximation of $u$ is a very challenging task for the reasons already mentioned above, but the fact that (1.1a) has to be solved numerically on a time interval of length of $\mathcal{O}(\varepsilon^{-1})$ makes things even worse[1]. An additional difficulty is the fact that in the term $\frac{1}{\varepsilon}Eu$ in (1.1a), which causes oscillations in time, $E$ is not a scalar, but a matrix with several different eigenvalues. In general, the operators $\mathcal{A}$ and $\frac{1}{\varepsilon}E$ do not commute, which means that some of the techniques used, e.g., in [2, 6, 11], are not suitable for our problem.

To tackle these challenges we combine *analytical* and *numerical* approximations. We cope with the spatial oscillations with the well-known *slowly varying envelope approximation* (SVEA) defined by

$$\widetilde{u}(t,x) = \mathrm{e}^{\mathrm{i}(\kappa \cdot x - \omega t)/\varepsilon}p(t,x) + \mathrm{e}^{-\mathrm{i}(\kappa \cdot x - \omega t)/\varepsilon}\overline{p(t,x)}, \tag{1.2}$$

[4, 14, 29, 32]. The number $\omega \in \mathbb{R}$ is chosen in such a way that the pair $(\omega, \kappa)$ fulfills the dispersion relation (cf. Section 2.1), and $p$ is the solution of a PDE called the *envelope equation*; cf. (2.5) below. Under a number of assumptions it was recently shown in [3] and [4] that the SVEA (1.2) approximates the solution $u$ of the original problem (1.1) up to an error of $\mathcal{O}(\varepsilon^2)$ in $L^\infty(\mathbb{R}^d)$ and on long time intervals $[0, t_{\mathrm{end}}/\varepsilon]$. The advantage is that $p$ is free from $\varepsilon$-induced oscillations in *space*, such that the space discretization of the envelope equation can be carried out with an $\varepsilon$-independent number of grid points. In contrast, time discretization is still a highly nontrivial problem, because the solution $p$ oscillates in time with the same frequency as $u$ and has to be approximated on the same long time interval $[0, t_{\mathrm{end}}/\varepsilon]$. Our approach is based on two steps. First, we replace the envelope equation by a new PDE, called the *reduced envelope equation* (REE), which has an additional advantage over the envelope equation when it comes to numerical approximation. After changing to co-moving coordinates, we devise a tailor-made, uniformly accurate integrator for the REE in the second step. Both steps require a careful analysis of interactions between oscillatory and non-oscillatory parts of the solution, which are identified by suitable projections. Substituting the numerical solution of the REE into the SVEA approximates the solution of the semilinear Friedrichs system up to an error of $\mathcal{O}(\tau + \varepsilon^2)$ in $L^\infty(\mathbb{R}^d)$, where $\tau$ is the step size. To the best of our knowledge, our approach is the first which achieves this accuracy with a number of time steps and grid points independent of the critical parameter $\varepsilon$. For *nonlinearly polarized* solutions (i.e. with initial data different from (1.1b)) higher-order approximations were recently constructed in [5].

---

[1]In Section 3.1 we will rescale time in such a way that computations only have to be made on $[0, t_{\mathrm{end}}]$. However, this rescaling does not make the problem easier, because in the new variables the frequency of the temporal oscillations is proportional to $\varepsilon^{-2}$ instead of $\varepsilon^{-1}$.

**Outline of the paper.** In the next section, we specify the problem setting and the analytical framework. Then, we introduce the envelope equation, discuss the accuracy of the SVEA, and we show that this accuracy does not deteriorate if the envelope equation is replaced by the REE. In Section 3, we construct an uniformly accurate time integrator for the REE, and we prove an error bound for the time discretization. Up to this point, the time integrator is stated in a compact but somewhat abstract form. In Section 4, we explain how this method can be turned into an implementable algorithm, which we apply to a one-dimensional version of the Maxwell–Lorentz system. The numerical examples corroborate the error bounds and illustrate the advantageous properties of our approach. Throughout the paper, we focus on the discretization in time, but in Section 4.2 we sketch how the space discretization can be carried out.

**Notation.** In the following the Euclidean scalar product of vectors $v_1, v_2 \in \mathbb{C}^m$ is represented by $v_1 \cdot v_2 = v_1^* v_2$. The $q$-norm on $\mathbb{C}^m$ or the induced matrix norm are both denoted by $|\cdot|_q$. The letter $I$ marks the identity matrix or identity operator, and i is the imaginary unit. For functions $f = f(t, x)$ that depend on time and space, we denote the mapping $x \mapsto f(t, x)$ by $f(t)$ rather than $f(t, \cdot)$. In a similar manner, for the spatial Fourier transform $\widehat{f}(t, k)$ of such a function, the dependence on $k$ is not explicitly specified. To keep the notation short, we write $L^1$ and $L^\infty$ instead of $L^1(\mathbb{R}^d, \mathbb{C}^m)$ and $L^\infty(\mathbb{R}^d, \mathbb{C}^m)$, respectively. All constants $C$ may depend on $t_{\text{end}}$, but not on $\varepsilon$, the step size $\tau$, nor on the number of time steps. Moreover, the values of these constants may vary from one step to the next.

# 2 Problem setting, slowly varying envelope approximation and reduced envelope equation

## 2.1 Polarization condition and envelope equation

As mentioned in the introduction, the function $p$ in (1.2) is the solution of the envelope equation. In order to formulate this PDE, we need some preparation. For $\alpha \in \mathbb{R}$ and $\beta \in \mathbb{R}^d$ we define the Hermitian matrix

$$L(\alpha, \beta) = -\alpha I + A(\beta) - iE \in \mathbb{C}^{m \times m} \qquad \text{with} \qquad A(\beta) = \sum_{j=1}^{d} A_j \beta_j. \tag{2.1}$$

From now on, we fix the wave vector $\kappa \in \mathbb{R}^d \setminus \{0\}$, which appears in (1.1b), and we choose $\omega$ to be an eigenvalue of $A(\kappa) - iE$. Hence, the kernel of $L(\omega, \kappa)$ is nontrivial, and the pair $(\omega, \kappa)$ is said to fulfill the dispersion relation.

**Assumption 2.1** (Polarization condition)
*The initial data in (1.1b) have the property that*

$$p_0(x) \in \ker\big(L(\omega, \kappa)\big) \text{ for almost all } x \in \mathbb{R}^d. \tag{2.2}$$

**Remark 2.2**
*This assumption was also made, e.g., in [3, 4, 14, 29]. Our results could also be extended to initial data of the form $p_0 = p_0^{(0)} + \varepsilon p_0^{(1)}$ with $p_0^{(0)} \in \ker\big(L(\omega, \kappa)\big)$, but since $\varepsilon$-dependent initial data makes the formulation more involved, we restrict ourselves to (2.2).*

Henceforth, we consider the trilinear extension of the real nonlinearity $T : \mathbb{R}^m \times \mathbb{R}^m \times \mathbb{R}^m \to \mathbb{R}^m$ to $\mathbb{C}^m \times \mathbb{C}^m \times \mathbb{C}^m$ and denote this extension again by $T$. Moreover, we define the symmetrized nonlinearity

$$T^{\text{sym}} : \mathbb{C}^m \times \mathbb{C}^m \times \mathbb{C}^m \to \mathbb{C}^m, \tag{2.3a}$$

$$T^{\text{sym}}(f_1, f_2, f_3) := T(f_1, f_2, \overline{f_3}) + T(f_1, \overline{f_2}, f_3) + T(\overline{f_1}, f_2, f_3). \tag{2.3b}$$

If $f_1 = f_2 = f_3 = f$, then we write $T^{\mathrm{sym}}(f)$ instead of $T^{\mathrm{sym}}(f, f, f)$. By definition, it follows that

$$T^{\mathrm{sym}}(g_1 + g_2, f_2, f_3) = T^{\mathrm{sym}}(g_1, f_2, f_3) + T^{\mathrm{sym}}(g_2, f_2, f_3), \tag{2.4a}$$

$$T^{\mathrm{sym}}(f_1, g_1 + g_2, f_3) = T^{\mathrm{sym}}(f_1, g_1, f_3) + T^{\mathrm{sym}}(f_1, g_2, f_3), \tag{2.4b}$$

$$T^{\mathrm{sym}}(f_1, f_2, g_1 + g_2) = T^{\mathrm{sym}}(f_1, f_2, g_1) + T^{\mathrm{sym}}(f_1, f_2, g_2). \tag{2.4c}$$

Since (2.3b) involves complex conjugation of different arguments, however, $T^{\mathrm{sym}}$ is only *real*-trilinear: for every *real* $c \in \mathbb{R}$ we have

$$T^{\mathrm{sym}}(cf_1, f_2, f_3) = T^{\mathrm{sym}}(f_1, cf_2, f_3) = T^{\mathrm{sym}}(f_1, f_2, cf_3) = c\,T^{\mathrm{sym}}(f_1, f_2, f_3)$$

for all $f_1, f_2, f_3 \in \mathbb{C}^m$, but this is *not* true if $c \in \mathbb{C} \setminus \mathbb{R}$ is complex.

Substituting the SVEA (1.2) into (1.1) and discarding higher harmonics yields the *envelope equation*

$$\partial_t p + \frac{\mathrm{i}}{\varepsilon} L(\omega, \kappa) p + \mathcal{A}p = \varepsilon T^{\mathrm{sym}}(p), \qquad t \in (0, t_{\mathrm{end}}/\varepsilon], \; x \in \mathbb{R}^d \tag{2.5a}$$

with initial condition

$$p(0, x) = p_0(x), \tag{2.5b}$$

cf. [4, 14, 29, 32]. In contrast to (1.1b), the highly oscillatory phase $\mathrm{e}^{\mathrm{i}\kappa \cdot x/\varepsilon}$ has been eliminated from the initial data (2.5b), and thus the solution of the envelope equation (2.5a) does not oscillate in space any more. But since the term $\frac{\mathrm{i}}{\varepsilon} L(\omega, \kappa) p$ causes rapid oscillations in time, solving (2.5) with standard methods would still require a tiny step size and hence unacceptable numerical costs.

## 2.2 Analytical setting

Recall that the Wiener algebra of vector-valued functions is the space

$$W = \left\{ f \in \left( \mathcal{S}'(\mathbb{R}^d) \right)^m : \widehat{f} \in L^1 \right\}, \qquad \|f\|_W = \|\widehat{f}\|_{L^1} = \int_{\mathbb{R}^d} |\widehat{f}(k)|_2 \; \mathrm{d}k, \tag{2.6}$$

where $\widehat{f} = \mathcal{F}f$ denotes the Fourier transform

$$(\mathcal{F}f)(k) = (2\pi)^{-d/2} \int_{\mathbb{R}^d} f(x) \mathrm{e}^{-\mathrm{i}k \cdot x} \, \mathrm{d}x$$

of $f$. For every $s \in \mathbb{N}_0$ the space

$$W^s = \{ f \in W : \partial^\alpha f \in W \text{ for all } \alpha \in \mathbb{N}_0^d, |\alpha|_1 \leq s \},$$

$$\|f\|_{W^s} = \sum_{|\alpha|_1 \leq s} \|\partial^\alpha f\|_W.$$

is a Banach algebra with continuous embedding $W \hookrightarrow L^\infty$, cf. [14, Proposition 1] and [29, Proposition 3.2].

For every $c_1, c_2 \in \mathbb{R}$ the operator

$$\mathrm{i}c_1 L(\omega, \kappa) + c_2 \mathcal{A} \colon W^1 \to W$$

generates a strongly continuous group $(\exp(t[\mathrm{i}c_1 L(\omega, \kappa) + c_2 \mathcal{A}]))_{t \in \mathbb{R}}$ on $W$ defined by

$$\mathcal{F}\Big( \exp\Big( t[\mathrm{i}c_1 L(\omega, \kappa) + c_2 \mathcal{A}] \Big) f \Big)(k) = \exp\Big( t\mathrm{i}[c_1 L(\omega, \kappa) + c_2 A(k)] \Big) \widehat{f}(k), \tag{2.7}$$

where $A(k)$ is the matrix defined in (2.1). Note that $A(\mathrm{i}k) = \mathrm{i}A(k)$. The $\exp(\dots)$ on the right-hand side is the standard matrix exponential function. The fact that $[c_1 L(\omega, \kappa) + c_2 A(k)]$ is Hermitian for every $k \in \mathbb{R}^d$ implies that

$$\left| \exp\Big( t\mathrm{i}[c_1 L(\omega, \kappa) + c_2 A(k)] \Big) \right|_2 = 1$$

and hence that

$$\left\| \exp \left( t[\mathrm{i}c_1 L(\omega, \kappa) + c_2 \mathcal{A}] \right) \right\|_{W^s} = 1, \qquad s \in \mathbb{N}_0. \tag{2.8}$$

Since $T$ is trilinear and $W$ is an algebra, there is a constant $C_T$ such that

$$\|T(f_1, f_2, f_3)\|_W \leq C_T \prod_{i=1}^{3} \|f_i\|_W \qquad \text{for all } f_1, f_2, f_3 \in W. \tag{2.9}$$

For $j \in \{1, \ldots, d\}$ we introduce the Fourier multiplier $(\mathcal{D}_j \widehat{f})(k) = \mathrm{i}k_j \widehat{f}(k)$, such that by definition $\mathcal{F}(\partial_j f) = \mathcal{D}_j \widehat{f}$ and, e.g.,

$$\|f\|_{W^1} = \|\widehat{f}\|_{L^1} + \sum_{j=1}^{d} \|\mathcal{D}_j \widehat{f}\|_{L^1}. \tag{2.10}$$

Local wellposedness in $W^s$ of the Friedrichs system (1.1) on long time intervals $[0, t_{\mathrm{end}}/\varepsilon]$ can be shown with standard arguments.

The matrix $L(\alpha, \beta)$ defined in (2.1) is of special interest, because it appears in the envelope equation (2.5) and in the polarization condition (Assumption 2.1). From now on, we assume the following.

**Assumption 2.3**

(i) The matrix $L(3\omega, 3\kappa)$ is regular and has no common eigenvalues with $L(\omega, \kappa)$.

(ii) The matrix $L(0, \beta)$ has a smooth eigendecomposition in the sense that the eigenvalues $\omega_\ell(\beta)$ and the corresponding eigenvectors $\phi_\ell(\beta)$ can be chosen in such a way that

$$\omega_\ell \in C^\infty(\mathbb{R}^d \setminus \{0\}, \mathbb{R}) \quad \text{and} \quad \phi_\ell \in C^\infty(\mathbb{R}^d \setminus \{0\}, \mathbb{C}^m) \qquad \text{for all } \ell = 1, \ldots, m.$$

With no loss of generality, we assume that for every $\beta$ the eigenvectors are pairwise orthogonal and normalized with respect to the Euclidean norm $|\cdot|_2$.

(iii) The eigenvalues $\omega_\ell(\beta)$ of $L(0, \beta)$ are globally Lipschitz continuous: there is a constant $C$ such that

$$|\omega_\ell(\tilde{\beta}) - \omega_\ell(\beta)| \leq C|\tilde{\beta} - \beta|_1 \qquad \text{for all } \tilde{\beta}, \beta \in \mathbb{R}^d.$$

(iv) Let $m_0 \in \{1, \ldots, m-1\}$ be the dimension of $\ker(L(\omega, \kappa))$. We always choose the enumeration of the eigenvalues in such a way that

$$\omega = \omega_\ell(\kappa) \quad \text{for } \ell = 1, \ldots, m_0. \tag{2.11}$$

(v) The eigenvalue $\omega = \omega_1(\kappa) = \ldots = \omega_{m_0}(\kappa)$ is bounded away from the other eigenvalues: There is a constant $C$ such that

$$|\omega - \omega_\ell(\beta)| \geq C \qquad \text{for all } \beta \in \mathbb{R}^d \text{ and } \ell = m_0 + 1, \ldots, m.$$

These assumptions are fulfilled, e.g., in case of the Maxwell–Lorentz system if $\omega$ is the largest or smallest eigenvalue; cf. [14, Example 3] and [4, Remark 2.4].

For fixed $\omega, \kappa$ and a vector $\theta \in \mathbb{R}^d$ we denote the eigenvalues of $L(\omega, \kappa + \theta)$ by $\lambda_\ell(\theta)$ and the associated eigenvectors by $\psi_\ell(\theta)$. These eigenvalues and eigenvectors will be used in Sections 2.3 and 4. Since $L(\omega, \kappa + \theta) = -\omega I + L(0, \kappa + \theta)$ by definition, it follows that $\lambda_\ell(\theta) = \omega_\ell(\kappa + \theta) - \omega$ and $\psi_\ell(\theta) = \phi_\ell(\kappa + \theta)$. Moreover, (2.11) implies that $\lambda_\ell(0) = 0$ for $\ell = 1, \ldots, m_0$, and that $\{\psi_1(0), \ldots, \psi_{m_0}(0)\}$ is an orthonormal basis of $\ker(L(\omega, \kappa))$.

## 2.3 Frequency-dependent projection and accuracy of the slowly varying envelope approximation

Our approach to solving the envelope equation (2.5) is based on the analysis carried out in [3, 4]. For the convenience of the reader, we briefly summarize the main results. Henceforth, we will always assume that the initial data $p_0$ are in $W^s$ with $s = 1$ or $s = 2$. With the classical fixed-point argument, it can be shown that there is a $t_{\text{end}} > 0$ such that a unique and uniformly bounded classical solution

$$p \in C([0, t_{\text{end}}/\varepsilon], W^1) \cap C^1([0, t_{\text{end}}/\varepsilon], W), \tag{2.12a}$$

$$\sup_{t \in [0, t_{\text{end}}/\varepsilon]} \|p(t)\|_{W^1} \leq C \quad \text{for all } \varepsilon \in (0, 1] \tag{2.12b}$$

of (2.5) with $C$ independent of $\varepsilon$ exists. If $p_0 \in W^2$, then we even have

$$p \in C([0, t_{\text{end}}/\varepsilon], W^2) \cap C^1([0, t_{\text{end}}/\varepsilon], W^1) \cap C^2([0, t_{\text{end}}/\varepsilon], W) \tag{2.13a}$$

$$\sup_{t \in [0, t_{\text{end}}/\varepsilon]} \|p(t)\|_{W^2} \leq C \quad \text{for all } \varepsilon \in (0, 1], \tag{2.13b}$$

cf. [14, Lemma 1], [3, Lemma 3.6.1], and [4, Lemma 2.2].

Before we discuss the properties of the SVEA and the envelope equation, we have to introduce some notation. Let[2]

$$\mathcal{P}_\varepsilon : L^1 \to L^1, \qquad \widehat{f}(k) \mapsto \sum_{\ell=1}^{m_0} \psi_\ell(\varepsilon k)\psi_\ell^*(\varepsilon k)\widehat{f}(k) \tag{2.14}$$

be the pointwise projection of $\widehat{f} \in L^1$ in Fourier space onto the first eigenspace of $L(\omega, \kappa + \varepsilon k)$ and denote the projector onto the orthogonal complement by $\mathcal{P}_\varepsilon^\perp = I - \mathcal{P}_\varepsilon$. In [14, Lemma 3] and [4, Section 3] the following bounds were shown.

**Proposition 2.4**
*Let $p$ be the solution of (2.5) with initial data $p_0$, and suppose that Assumptions 2.1 and 2.3 hold.*

(i) *If $p_0 \in W^1$, then there is a constant $C > 0$ such that*

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|\mathcal{P}_\varepsilon^\perp \widehat{p}(t)\|_{L^1} \leq C\varepsilon \tag{2.15}$$

*for all $\varepsilon \in (0, 1]$.*

(ii) *If in addition $p_0 \in W^2$, then there is a constant $C > 0$ such that*

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|\mathcal{D}_j \mathcal{P}_\varepsilon^\perp \widehat{p}(t)\|_{L^1} \leq C\varepsilon \tag{2.16}$$

*for all $\varepsilon \in (0, 1]$ and every $j \in \{1, \ldots, d\}$.*

According to (2.10) the bounds (2.15) and (2.16) imply that $\|\mathcal{F}^{-1}(\mathcal{P}_\varepsilon^\perp \widehat{p}(t))\|_{W^1} = \mathcal{O}(\varepsilon)$, which means that the projected part $\mathcal{P}_\varepsilon^\perp \widehat{p}$ of the solution of (2.5) remains small even on long time intervals. This property plays a crucial role in our approach.

With Proposition 2.4 the following error bound for the SVEA was proved in [4, Theorem 4.3] and similar in [3, Theorem 4.3.4].

**Theorem 2.5** (Error bound for the SVEA)
*Let $p_0 \in W^2$ and let $u$ be the solution of (1.1). Let $p$ be the solution of the envelope equation (2.5) and let $\widetilde{u}$ be the approximation defined in (1.2). Under Assumptions 2.1 and 2.3 there is a constant $C$ such that*

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|u(t) - \widetilde{u}(t)\|_W \leq C\varepsilon^2, \tag{2.17}$$

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|u(t) - \widetilde{u}(t)\|_{L^\infty} \leq C\varepsilon^2. \tag{2.18}$$

---

[2]We write $\widehat{f}$ instead of $f$, because later the function will be the Fourier transform of $f \in W$.

Theorem 2.5 suggests to approximate the solution of the Friedrichs system (1.1) by solving the envelope equation (2.5) numerically and then replacing $p$ in the SVEA (1.2) by its numerical approximation. Then, the total error consists of the *a priori* error of $\mathcal{O}(\varepsilon^2)$ from (2.18) plus the numerical error caused by approximating $p$ numerically. This is a viable approach, but it turns out that for the construction of a numerical method it is more advantageous to replace the envelope equation by yet another PDE which is more suitable for time integration, and which yields the same *a priori* error; cf. Remark 4.1 below. This PDE, called the *reduced envelope equation*, is derived in the next subsection.

## 2.4 Frequency-independent projection and reduced envelope equation

With the projection $\mathcal{P}_\varepsilon$ the Fourier transform $\widehat{p}$ of the solution $p$ of the envelope equation can be decomposed into the two parts $\mathcal{P}_\varepsilon\widehat{p}$ and $\mathcal{P}_\varepsilon^\perp\widehat{p}$. It was shown in [14, Lemma 2] and [4, Lemma 3.5] that $\mathcal{P}_\varepsilon\widehat{p}$ is "essentially non-oscillatory" in the sense that the first two time derivatives are uniformly bounded in $\varepsilon$. Hence, the oscillatory behavior of $p$ comes only from the part $\mathcal{P}_\varepsilon^\perp\widehat{p}$, which, however, is small by Proposition 2.4. Our ansatz for the derivation of the new PDE is, roughly speaking, to omit as many oscillatory but small parts as possible in the nonlinearity at the cost of an error of $\mathcal{O}(\varepsilon^2)$. After denoting the "smooth" and the oscillatory parts of $p$ by $p_{\mathrm{smo}} = \mathcal{F}^{-1}(\mathcal{P}_\varepsilon\widehat{p})$ and $p_{\mathrm{osc}} = \mathcal{F}^{-1}(\mathcal{P}_\varepsilon^\perp\widehat{p})$, respectively, we infer from (2.4) and (2.9) that formally

$$
\begin{aligned}
T^{\mathrm{sym}}(p) = T^{\mathrm{sym}}(p,p,p) &= T^{\mathrm{sym}}\Big(p_{\mathrm{smo}} + p_{\mathrm{osc}},\ p_{\mathrm{smo}} + p_{\mathrm{osc}},\ p_{\mathrm{smo}} + p_{\mathrm{osc}}\Big) \\
&= T^{\mathrm{sym}}\Big(p_{\mathrm{smo}}, p_{\mathrm{smo}}, p_{\mathrm{smo}}\Big) + T^{\mathrm{sym}}\Big(p_{\mathrm{smo}}, p_{\mathrm{smo}}, p_{\mathrm{osc}}\Big) \\
&\quad + T^{\mathrm{sym}}\Big(p_{\mathrm{smo}}, p_{\mathrm{osc}}, p_{\mathrm{smo}}\Big) + T^{\mathrm{sym}}\Big(p_{\mathrm{osc}}, p_{\mathrm{smo}}, p_{\mathrm{smo}}\Big) + \mathcal{O}(\varepsilon^2).
\end{aligned} \tag{2.19}
$$

In this representation, the oscillatory part $p_{\mathrm{osc}}$ appears in at most one of the three components of $T^{\mathrm{sym}}$. This is an advantage for the time integration, because approximating oscillatory functions is more intricate than smooth ones. The problem is that the nonlinearity $T^{\mathrm{sym}}$ is defined for functions in physical space, whereas the computation of $p_{\mathrm{smo}}$ and $p_{\mathrm{osc}}$ requires a forward and inverse Fourier transform, because the projection $\mathcal{P}_\varepsilon$ can only be applied in Fourier space. Unfortunately, this "detour" is not compatible with the techniques used in the construction of the time integration as presented in Section 3.2 below. For this reason, we will now consider a similar but frequency-independent projection which has essentially the same beneficial properties.

Since $\{\psi_1(0),\ \ldots,\ \psi_{m_0}(0)\}$ is an orthonormal basis of $\ker(L(\omega,\kappa))$, it follows that

$$
v \mapsto Pv, \qquad P = \sum_{\ell=1}^{m_0} \psi_\ell(0)\psi_\ell^*(0) \in \mathbb{C}^{m\times m} \tag{2.20}
$$

is the orthogonal projection from $\mathbb{C}^m$ onto the kernel of $L(\omega,\kappa)$. The essential difference between the projections (2.14) and (2.20) is that $P$ does not depend on $\varepsilon$ nor on the frequency $k$, but only on the matrix $L(\omega,\kappa)$. Hence, (2.20) can be applied to an arbitrary *vector* $v \in \mathbb{C}^m$, whereas (2.14) is only defined for a vector-valued *function* in Fourier space. Of course, $P$ can also be extended to an operator

$$
P : L^1 \to L^1, \qquad \widehat{f}(k) \mapsto \sum_{\ell=1}^{m_0} \psi_\ell(0)\psi_\ell^*(0)\widehat{f}(k) \quad \text{for all } k \in \mathbb{R}^d, \tag{2.21}
$$

which maps functions to functions. Although strictly speaking (2.20) and (2.21) are two different mappings, we will denote both with the same symbol $P$. Below, we will often use that $\|Pf\|_W \le \|f\|_W$ and $\|P^\perp f\|_W \le \|f\|_W$, which follows by definition.

Our goal is to show that Proposition 2.4 remains true if the frequency-dependent projection $\mathcal{P}_\varepsilon$ is replaced by $P$. As a first step, we quote the following estimate for the difference between $\mathcal{P}_\varepsilon$ and $P$.

**Lemma 2.6**
*Under Assumption 2.3, there is a constant $C$ such that the bound*

$$
\|\mathcal{P}_\varepsilon\widehat{f} - P\widehat{f}\|_{L^1} \le C\varepsilon\|f\|_{W^1}
$$

*holds for all $f \in W^1$ and all $\varepsilon \in (0,1]$.*

This result was shown as a part of the proof of Lemma 3 in [14].

With Lemma 2.6 we can now prove the following counterpart of Proposition 2.4.

**Corollary 2.7**
*Let $p$ be the solution of (2.5). Under the assumptions of Proposition 2.4(ii) there is a constant $C$ such that*

$$\sup_{t\in[0,t_{end}/\varepsilon]} \left\| P^\perp p(t) \right\|_{W^1} \leq C\varepsilon$$

*for all $\varepsilon \in (0,1]$.*

*Proof.* With $I - P = I - \mathcal{P}_\varepsilon + (\mathcal{P}_\varepsilon - P)$, Proposition 2.4(i) and Lemma 2.6 imply that

$$\sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| P^\perp \widehat{p}(t) \right\|_{L^1} \leq \sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| \mathcal{P}_\varepsilon^\perp \widehat{p}(t) \right\|_{L^1} + \sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| (\mathcal{P}_\varepsilon - P) \widehat{p}(t) \right\|_{L^1}$$

$$\leq C\varepsilon + C\varepsilon \sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| p(t) \right\|_{W^1}.$$

For every $j = 1, \dots, d$, we repeat the calculation and use Proposition 2.4(ii) to obtain

$$\sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| \mathcal{D}_j P^\perp \widehat{p}(t) \right\|_{L^1} \leq C\varepsilon + C\varepsilon \sup_{t\in[0,t_{\mathrm{end}}/\varepsilon]} \left\| p(t) \right\|_{W^2}.$$

Now the assertion follows from the uniform boundedness of $\left\| p(t) \right\|_{W^2}$ and from the fact that $\left\| P^\perp p(t) \right\|_{W^1} = \left\| P^\perp \widehat{p}(t) \right\|_{L^1} + \sum_{j=1}^d \left\| \mathcal{D}_j P^\perp \widehat{p}(t) \right\|_{L^1}$. $\qquad\square$

**Lemma 2.8**
*Let $p_0 \in W^2$ and let $p$ be the solution of the envelope equation (2.5). Then, the time derivative of $Pp$ is uniformly bounded: there is a constant $C$ independent of $\varepsilon \in (0,1]$ such that*

$$\sup_{t\in[0,t_{end}/\varepsilon]} \left\| \partial_t Pp(t) \right\|_{W^1} \leq C.$$

We omit the proof, because the result can be shown in a similar way as Lemma 2.11 below.

Corollary 2.7 and Lemma 2.8 allow us to split the solution $p = Pp + P^\perp p$ of the envelope equation into an essentially non-oscillatory part $Pp$ and an oscillatory but small part $P^\perp p$. Hence, this decomposition has the same favorable properties as $p = p_{\mathrm{smo}} + p_{\mathrm{osc}}$ at the beginning of this subsection, but the computational disadvantages of the latter are now avoided, because no Fourier transform is required.

We are now in a position to replace the envelope equation (2.5a) by a new PDE which is even better suited for numerical computations, and which does not spoil the accuracy of the SVEA (cf. Theorem 2.5). For this purpose, we define $T_{\mathrm{perm}}^{\mathrm{sym}}(f_1, f_2, f_3)$ to be the sum of $T^{\mathrm{sym}}$ evaluated for all permutations of the arguments, i.e.

$$T_{\mathrm{perm}}^{\mathrm{sym}}(f_1, f_2, f_3) = T^{\mathrm{sym}}(f_1, f_2, f_3) + T^{\mathrm{sym}}(f_1, f_3, f_2) + T^{\mathrm{sym}}(f_2, f_1, f_3) \tag{2.22}$$
$$+ T^{\mathrm{sym}}(f_2, f_3, f_1) + T^{\mathrm{sym}}(f_3, f_1, f_2) + T^{\mathrm{sym}}(f_3, f_2, f_1).$$

With this notation and the decomposition $p = Pp + P^\perp p$, we can represent the symmetrized nonlinearity in the form

$$T^{\mathrm{sym}}(p) = T^{\mathrm{sym}}(Pp + P^\perp p)$$
$$= T^{\mathrm{sym}}(Pp) + \frac{1}{2} T_{\mathrm{perm}}^{\mathrm{sym}}(P^\perp p, Pp, Pp) + \frac{1}{2} T_{\mathrm{perm}}^{\mathrm{sym}}(P^\perp p, P^\perp p, Pp) + T^{\mathrm{sym}}(P^\perp p). \tag{2.23}$$

Since $P^\perp p = \mathcal{O}(\varepsilon)$ in the sense of Corollary 2.7, we neglect all terms where $P^\perp p$ appears in at least two of the three arguments, as in (2.19). This yields the *reduced envelope equation (REE)*

$$\partial_t q + \frac{\mathrm{i}}{\varepsilon} L(\omega, \kappa) q + \mathcal{A}q = \varepsilon T^{\mathrm{sym}}(Pq) + \frac{\varepsilon}{2} T_{\mathrm{perm}}^{\mathrm{sym}}(P^\perp q, Pq, Pq), \qquad t \in (0, t_{\mathrm{end}}/\varepsilon], x \in \mathbb{R}^d, \tag{2.24a}$$

$$q(0, x) = p_0(x). \tag{2.24b}$$

It can be shown that (2.12) and (2.13) hold true if $p$ is replaced by the solution $q$ of (2.24).

Our next goal is to prove that $q$ approximates $p$ up to an error of $\mathcal{O}(\varepsilon^2)$ in $\|\cdot\|_{W^1}$. This means, in particular, that Theorem 2.5 remains true if we replace $p$ by $q$ in the SVEA (1.2). As a preparation, we state bounds for the nonlinearities $T^{\mathrm{sym}}$ and $T^{\mathrm{sym}}_{\mathrm{perm}}$. If $r \in \mathbb{N}_0$ and $f_1, f_2, f_3, g_1, g_2, g_3 \in W^r$ with $\widehat{C} = \max\{\|f_1\|_{W^r}, \ldots, \|g_3\|_{W^r}\}$, then (2.9) and (2.4) imply the inequalities

$$\|T^{\mathrm{sym}}(f_1, f_2, f_3)\|_{W^r} \leq C \prod_{i=1}^{3} \|f_i\|_{W^r}, \tag{2.25}$$

$$\|T^{\mathrm{sym}}(f_1, f_2, f_3) - T^{\mathrm{sym}}(g_1, g_2, g_3)\|_{W^r} \leq C \sum_{i=1}^{3} \|f_i - g_i\|_{W^r} \tag{2.26}$$

for the symmetrized nonlinearity, as well as

$$\left\|T^{\mathrm{sym}}_{\mathrm{perm}}(f_1, f_2, f_3)\right\|_{W^r} \leq C \prod_{i=1}^{3} \|f_i\|_{W^r}, \tag{2.27}$$

$$\left\|T^{\mathrm{sym}}_{\mathrm{perm}}(f_1, f_2, f_3) - T^{\mathrm{sym}}_{\mathrm{perm}}(g_1, g_2, g_3)\right\|_{W^r} \leq C \sum_{i=1}^{3} \|f_i - g_i\|_{W^r}, \tag{2.28}$$

where $C$ depends on $C_T$ and $\widehat{C}$.

For the proof of the following theorem, we define the abbreviation.

$$\mathcal{E}_{\mathcal{A}}(t) := \exp\left(-t(\mathrm{i}L(\omega, \kappa) + \varepsilon\mathcal{A})\right). \tag{2.29}$$

**Theorem 2.9**
*Let $p_0 \in W^2$ and suppose that Assumptions 2.1 and 2.3 hold. Then, the difference between the solutions of (2.5) and (2.24) is bounded by*

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|p(t) - q(t)\|_{W^1} \leq C\varepsilon^2. \tag{2.30}$$

*Proof.* Let $C_{p,q}$ be a constant such that

$$\sup_{t \in [0, t_{\mathrm{end}}/\varepsilon]} \|p(t)\|_{W^1} \leq C_{p,q}, \qquad \sup_{t \in [0, t_{\mathrm{end}}/\varepsilon]} \|q(t)\|_{W^1} \leq C_{p,q} \qquad \text{for all } \varepsilon \in (0, 1]. \tag{2.31}$$

For every $t \in [0, t_{\mathrm{end}}/\varepsilon]$ Duhamel's formula yields the representations

$$p(t) = \mathcal{E}_{\mathcal{A}}\left(\tfrac{t}{\varepsilon}\right) p_0 + \varepsilon \int_0^t \mathcal{E}_{\mathcal{A}}\left(\tfrac{t-s}{\varepsilon}\right) T^{\mathrm{sym}}(p(s)) \, \mathrm{d}s,$$

$$q(t) = \mathcal{E}_{\mathcal{A}}\left(\tfrac{t}{\varepsilon}\right) p_0 + \varepsilon \int_0^t \mathcal{E}_{\mathcal{A}}\left(\tfrac{t-s}{\varepsilon}\right) \left(T^{\mathrm{sym}}(Pq(s)) + \frac{1}{2} T^{\mathrm{sym}}_{\mathrm{perm}}\left(P^\perp q(s), Pq(s), Pq(s)\right)\right) \mathrm{d}s.$$

Since $\mathcal{E}_{\mathcal{A}}\left(\tfrac{t}{\varepsilon}\right)$ is an isometry on $W^1$ for all $t \in \mathbb{R}$ (cf. (2.8)), it follows from (2.23) that

$$\|p(t) - q(t)\|_{W^1}$$
$$\leq \varepsilon \int_0^t \left\|T^{\mathrm{sym}}(p(s)) - T^{\mathrm{sym}}(Pq(s)) - \frac{1}{2} T^{\mathrm{sym}}_{\mathrm{perm}}\left(P^\perp q(s), Pq(s), Pq(s)\right)\right\|_{W^1} \mathrm{d}s$$
$$\leq \varepsilon \int_0^t \|T^{\mathrm{sym}}(Pp(s)) - T^{\mathrm{sym}}(Pq(s))\|_{W^1} \, \mathrm{d}s$$

$$+ \frac{\varepsilon}{2} \int\limits_0^t \left\| T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp p(s), Pp(s), Pp(s)) - T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(s), Pq(s), Pq(s)) \right\|_{W^1} \mathrm{d}s$$

$$+ \varepsilon \int\limits_0^t \left\| \frac{1}{2} T^{\mathrm{sym}}_{\mathrm{perm}} \left(P^\perp p(s), P^\perp p(s), Pp(s)\right) + T^{\mathrm{sym}}\left(P^\perp p(s)\right) \right\|_{W^1} \mathrm{d}s.$$

For the first term on the right-hand side, (2.26) and (2.31) imply the bound

$$\|T^{\mathrm{sym}}(Pp(s)) - T^{\mathrm{sym}}(Pq(s))\|_{W^1} \le C_1 \|p(s) - q(s)\|_{W^1}$$

with a constant $C_1$. For the second term the inequality

$$\left\| T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp p(s), Pp(s), Pp(s)) - T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(s), Pq(s), Pq(s)) \right\|_{W^1} \le C_2 \|p(s) - q(s)\|_{W^1}$$

follows from (2.28). With (2.27), (2.25), and Corollary 2.7, we can derive the estimate

$$\left\| \frac{1}{2} T^{\mathrm{sym}}_{\mathrm{perm}} \left(P^\perp p(s), P^\perp p(s), Pp(s)\right) + T^{\mathrm{sym}} \left(P^\perp p(s)\right) \right\|_{W^1} \le C_3 \varepsilon^2$$

for the last term. By combining these bounds, we obtain the inequality

$$\|p(t) - q(t)\|_{W^1} \le C_3 \varepsilon^2 + \varepsilon(C_1 + C_2) \int\limits_0^t \|p(s) - q(s)\|_{W^1} \mathrm{d}s$$

for the error, and applying Gronwall's lemma and using that $\varepsilon t \le t_{\mathrm{end}}$ proves (2.30). □

In the remainder of this subsection, we will show that $q$ has the same properties as $p$. The following result is the counterpart of Corollary 2.7.

### Corollary 2.10
*Let $q$ be the solution of (2.24) with $p_0 \in W^2$. Under Assumption 2.1 and 2.3 there is a constant $C$ such that*

$$\sup_{t \in [0, t_{end}]} \left\| P^\perp q(t) \right\|_{W^1} \le C \varepsilon$$

*for all $\varepsilon \in (0, 1]$.*

*Proof.* Since

$$\left\| P^\perp q(t) \right\|_{W^1} \le \left\| P^\perp q(t) - P^\perp p(t) \right\|_{W^1} + \left\| P^\perp p(t) \right\|_{W^1}$$

for all $t \in [0, t_{\mathrm{end}}/\varepsilon]$, the assertion follows directly from Theorem 2.9 and Corollary 2.7. □

Finally, we prove the counterpart of Lemma 2.8

### Lemma 2.11
*Let $p_0 \in W^2$ and let $q$ be the solution of the REE (2.24). Then, the time derivative of $Pq$ is uniformly bounded: there is a constant $C$ independent of $\varepsilon \in (0, 1]$ such that*

$$\sup_{t \in [0, t_{end}/\varepsilon]} \|\partial_t Pq(t)\|_{W^1} \le C.$$

*Proof.* We apply $P$ on both sides of (2.24a) and use that $PL(\omega, \kappa) = 0$ by definition. For every $t \in [0, t_{\mathrm{end}}/\varepsilon]$, this yields

$$\|\partial_t Pq(t)\|_{W^1} = \|P\mathcal{A}q(t)\|_{W^1} + \varepsilon \|PT^{\mathrm{sym}}(q(t))\|_{W^1} + \frac{\varepsilon}{2} \left\| PT^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(t), Pq(t), Pq(t)) \right\|_{W^1}$$

$$\le C \|q(t)\|_{W^2} + C\varepsilon \|q(t)\|^3_{W^1}$$

according to (2.25) and (2.27). Now the assertion follows from the uniform boundedness of $q$ and the fact that $\varepsilon \le 1$. □

The results of this subsection reveal why solving the REE (2.24a) numerically instead of the envelope equation (2.5a) is more advantageous in spite of the seemingly more complicated nonlinear term. In the modified nonlinearity on the right-hand side of (2.24a), the first term $T^{\text{sym}}(Pq)$ has a bounded time derivative, and in the second term $T^{\text{sym}}_{\text{perm}}(P^\perp q, Pq, Pq)$ the fast part $P^\perp q$ of the solution appears *in only one component*. This means, that approximating the oscillatory behavior of $T^{\text{sym}}_{\text{perm}}(P^\perp q, Pq, Pq)$ is less involved than for the full nonlinearity $T^{\text{sym}}(p) = T^{\text{sym}}(p, p, p)$, where all three components oscillate with frequency of $\mathcal{O}(\varepsilon^{-1})$ and interact in a complicated way.

# 3 Time integration

In the previous section we have shown that the solution $u$ of the semilinear Friedrichs system (1.1) can be approximated by

$$\widetilde{u}_{\text{REE}}(t,x) = e^{i(\kappa \cdot x - \omega t)/\varepsilon} q(t,x) + e^{-i(\kappa \cdot x - \omega t)/\varepsilon} \overline{q(t,x)}, \tag{3.1}$$

where $q$ solves the REE (2.24). Theorems 2.5 and 2.9 yield the error bound

$$\sup_{t \in [0, t_{\text{end}}/\varepsilon]} \|u(t) - \widetilde{u}_{\text{REE}}(t)\|_{L^\infty} \leq C\varepsilon^2.$$

In this and the next section, we construct and analyze a novel time integrator which for $n \in \mathbb{N}$ computes approximations to $q(t_n, x)$ at $t_n = n\tau$ with a step size $\tau > 0$ and a numerical error of $\mathcal{O}(\tau)$ in $L^\infty(\mathbb{R}^d)$. Evaluating (3.1) for $t = t_n$ and replacing $q(t,x)$ by the numerical approximation leads to a total error of $\mathcal{O}(\varepsilon^2 + \tau)$ in $L^\infty(\mathbb{R}^d)$.

## 3.1 Change of variable

For numerical simulations we need to truncate the unbounded spatial domain $\mathbb{R}^d$ in a suitable way; cf. Section 4.2. The solution of the REE is a wave packet which propagates with group velocity $c_g = \nabla_\kappa \omega(\kappa)$, such that tracking the solution on the time interval $[0, t_{\text{end}}/\varepsilon]$ would require a huge computational domain of size of $\mathcal{O}(\varepsilon^{-1})$, roughly speaking. It is much more efficient to transform the REE to co-moving variables. Moreover, we rescale time in such a way that the evolution is considered on the interval $[0, t_{\text{end}}]$ instead of $[0, t_{\text{end}}/\varepsilon]$. The reason is that on an $\varepsilon$-dependent time interval the number of time steps does not only depend on the step size, but also on $\varepsilon$, and hence the step size alone does not give any information about the total numerical costs. Moreover, our goal is to construct a uniformly accurate integrator of order one for the REE, which means that the constant in the error bound must not depend on $\varepsilon$. But since such an error constant always depends on the length of the time interval, we have to work on the rescaled time interval $[0, t_{\text{end}}]$.

For these reasons, we use the change of variables

$$\mathbf{x} = x - tc_g, \qquad \mathbf{t} = \varepsilon t \in [0, t_{\text{end}}], \qquad \mathbf{q}(\mathbf{t}, \mathbf{x}) = q(t, x).$$

Substituting into the REE (2.24) yields that the new function $\mathbf{q}(\mathbf{t}, \mathbf{x})$ is the solution of

$$\partial_{\mathbf{t}} \mathbf{q} + \frac{i}{\varepsilon^2} L(\omega, \kappa) \mathbf{q} + \frac{1}{\varepsilon} \Big( \mathcal{A}_{\mathbf{x}} - c_g \cdot \nabla_{\mathbf{x}} \Big) \mathbf{q} = T^{\text{sym}}(P\mathbf{q}) + \frac{1}{2} T^{\text{sym}}_{\text{perm}}(P^\perp \mathbf{q}, P\mathbf{q}, P\mathbf{q}),$$
$$\mathbf{q}(0, \mathbf{x}) = p_0(x)$$

with $\mathcal{A}_{\mathbf{x}} = \sum_{j=1}^d A_j \frac{\partial}{\partial \mathbf{x}_j}$. However, since the function $q(t, x)$ in the original variables will not appear again in the remainder of this paper, we can immediately drop the boldface notation and write $t$, $x$, $q$, $\mathcal{A}$ instead of $\mathbf{t}$, $\mathbf{x}$, $\mathbf{q}$, $\mathcal{A}_{\mathbf{x}}$ for the sake of simplicity. Note that the choice of variables does not affect the initial data, because $(\mathbf{t}, \mathbf{x}) = (t, x)$ for $t = 0$. Moreover, we will henceforth omit the arguments of $L(\omega, \kappa)$, because the wave vector $\kappa \in \mathbb{R}^d \setminus \{0\}$ and the eigenvalue $\omega(\kappa) \in \mathbb{R}$ are kept fixed from now on. Finally, we define the new differential operator $\mathcal{B} = \mathcal{A} - c_g \cdot \nabla$. With these changes and conventions, the REE

takes the new form

$$\partial_t q + \frac{\mathrm{i}}{\varepsilon^2} L q + \frac{1}{\varepsilon} \mathcal{B} q = T^{\mathrm{sym}}(Pq) + \frac{1}{2} T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q, Pq, Pq), \quad t \in (0, t_{\mathrm{end}}], x \in \mathbb{R}^d \tag{3.2a}$$

$$q(0, x) = p_0(x). \tag{3.2b}$$

Analogous to (2.29) we set

$$\mathcal{E}_{\mathcal{B}}(t) := \exp\left(-t(\mathrm{i}L + \varepsilon \mathcal{B})\right). \tag{3.3}$$

## 3.2 Construction of the method

The aim is now to construct a uniformly accurate time integrator for the REE (3.2) to compute approximations $q^n \approx q(t_n)$ for times $t_n = n\tau$ with step size $\tau = t_{\mathrm{end}}/N$, where $N \in \mathbb{N}$ is the total number of time steps. We will always assume that

$$q \in C([0, t_{\mathrm{end}}], W^2) \cap C^1([0, t_{\mathrm{end}}], W^1) \cap C^2([0, t_{\mathrm{end}}], W)$$

and that there is a constant $C$ such that

$$\sup_{t \in [0, t_{\mathrm{end}}]} \|q(t)\|_{W^2} \leq C \quad \text{for all } \varepsilon \in (0, 1]. \tag{3.4}$$

It can be shown that this is true if $p_0 \in W^2$ and $t_{\mathrm{end}}$ is sufficiently small; cf. Section 2.4.

From now on, the notation

$$f(t, \varepsilon, \tau) = \mathcal{O}(\varepsilon^j \tau^\ell)$$

for an error term $f(t, \varepsilon, \tau)$ and $j, \ell \in \mathbb{N}_0$ means that there is a constant $C$ such that

$$\sup_{t \in [0, t_{\mathrm{end}}]} \|f(t, \varepsilon, \tau)\|_W \leq C \varepsilon^j \tau^\ell.$$

As a starting point for the construction we consider Duhamel's formula

$$q(t_n + \tau) = \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau}{\varepsilon^2}\right) q(t_n) + \int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau - s}{\varepsilon^2}\right) T^{\mathrm{sym}}(Pq(t_n + s)) \, \mathrm{d}s \tag{3.5}$$

$$+ \frac{1}{2} \int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau - s}{\varepsilon^2}\right) T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(t_n + s), Pq(t_n + s), Pq(t_n + s)) \, \mathrm{d}s.$$

We have to approximate the integrals up to an error of $\mathcal{O}(\tau^2)$ by computable expressions which involve only $q(t_n)$.

The smooth part $Pq(t_n + s)$ can simply be replaced by $Pq(t_n)$, because Lemma 2.11 implies that

$$\|Pq(t_n + s) - Pq(t_n)\|_{W^1} \leq Cs, \qquad s \geq 0 \tag{3.6}$$

if $p_0 \in W^2$. The uniform boundedness of $q$ and the inequalities (2.26) and (2.28) yield

$$q(t_n + \tau) = \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau}{\varepsilon^2}\right) q(t_n) + \int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau - s}{\varepsilon^2}\right) T^{\mathrm{sym}}(Pq(t_n)) \, \mathrm{d}s \tag{3.7}$$

$$+ \frac{1}{2} \int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau - s}{\varepsilon^2}\right) T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(t_n + s), Pq(t_n), Pq(t_n)) \, \mathrm{d}s + \mathcal{O}(\tau^2).$$

because as in (2.8) one can check that

$$\left\| \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \right\|_W = \left\| \exp\left(-\tfrac{(\tau-s)}{\varepsilon^2}(\mathrm{i}L + \varepsilon\mathcal{B})\right) \right\|_W = 1.$$

The first integral on the right-hand side of (3.7) can then be computed by means of the Fourier transform of $T^{\mathrm{sym}}(Pq(t_n))$; cf. Section 4.

The main challenge is to approximate the second integral term

$$\frac{1}{2}\int_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) g(s)\,\mathrm{d}s, \qquad g(s) = T^{\mathrm{sym}}_{\mathrm{perm}}(P^\perp q(t_n + s), Pq(t_n), Pq(t_n)) \tag{3.8}$$

accurately and efficiently. The action of the operator

$$\mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) = \exp\left(-\frac{(\tau-s)}{\varepsilon^2}(\mathrm{i}L + \varepsilon\mathcal{B})\right) \tag{3.9}$$

on a function has to be computed in Fourier space due to the differential operator $\mathcal{B}$; cf. (2.7). Hence, we have to approximate the Fourier transform $\widehat{g}(s)$ of $g(s)$ in some way or the other. The first option is to consider the counterpart of $T^{\mathrm{sym}}_{\mathrm{perm}}$ in Fourier space, i.e. to derive a nonlinear operator $\widetilde{T}^{\mathrm{sym}}_{\mathrm{perm}}$ with the property that

$$\mathcal{F}g(s) = \widetilde{T}^{\mathrm{sym}}_{\mathrm{perm}}\Big(\mathcal{F}P^\perp q(t_n + s), \mathcal{F}Pq(t_n), \mathcal{F}Pq(t_n)\Big).$$

But since the nonlinearitiy $T$ and thus also $T^{\mathrm{sym}}_{\mathrm{perm}}$ typically involve multiplications (as, e.g., in the Maxwell–Lorentz model (4.7)) and since the Fourier transform turns multiplications into convolutions, computing (3.8) by means of $\widetilde{T}^{\mathrm{sym}}_{\mathrm{perm}}$ would cause numerical costs which grow cubically with the number of grid points; cf. [3, Equations (3.23),(5.12) and Section 5.5.1]. For this reason, it is crucial to first compute $g(s)$ in physical space and only then transform the result to Fourier space.

For an approximation of the integral (3.8), it is tempting to freeze $g(s)$ at some constant value, say $g(s) \approx g(0)$, and solve the resulting integral analytically, as for the first integral on the right-hand side of (3.7). Unfortunately, freezing $g(s) \approx g(0)$ spoils the accuracy because $g$ involves the highly oscillatory part $P^\perp q(t_n + s)$. Hence, we need a better approximation to $P^\perp q(t_n + s)$ which satisfies the following two criteria: First, evaluations of the resulting counterpart of $g(s)$ must be cheap, and second, the corresponding integral must be analytically computable. In contrast to [6, 10, 11] and many other works, we cannot simply use Duhamel's formula (3.5) a second time for this purpose, because this would lead to expressions of the form

$$T^{\mathrm{sym}}_{\mathrm{perm}}\left(P^\perp \mathcal{E}_\mathcal{B}\left(\tfrac{s}{\varepsilon^2}\right) q(t_n), \ \ldots \ , \ \ldots \right).$$

Such terms *cannot* be computed efficiently, because $\mathcal{E}_\mathcal{B}\left(\tfrac{s}{\varepsilon^2}\right)$ cannot be evaluated in physical space, as we have pointed out before.

At this point, the crucial idea is to employ the representation

$$P^\perp q(t_n + s) = E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q(t_n) + \mathcal{I}_1^n(s, q, \varepsilon) + \mathcal{I}_2^n(s, q, \varepsilon) + \mathcal{I}_3^n(s, q, \varepsilon) \tag{3.10}$$

with[3]

$$E(t) = \exp(-\mathrm{i}tL) \qquad \text{for } t \in \mathbb{R}, \tag{3.11}$$

$$\mathcal{I}_1^n(s, q, \varepsilon) = -\frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) P^\perp \mathcal{B}q(t_n + r)\,\mathrm{d}r,$$

$$\mathcal{I}_2^n(s, q, \varepsilon) = \int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) P^\perp T^{\mathrm{sym}}(Pq(t_n + r))\,\mathrm{d}r,$$

$$\mathcal{I}_3^n(s, q, \varepsilon) = \frac{1}{2}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) P^\perp T^{\mathrm{sym}}_{\mathrm{perm}}\big(P^\perp q(t_n + r), Pq(t_n + r), Pq(t_n + r)\big)\,\mathrm{d}r.$$

---

[3]The time-dependent matrix $E(t)$ is, of course, not the same as the matrix $E$ in (1.1a).

This is again obtained by applying Duhamel's formula to (3.2a), but with the term $\frac{1}{\varepsilon}\mathcal{B}q$ removed from the linear part and grouped together with the nonlinear terms instead. The advantage of the representation (3.10) is that $E(t)$ is only a matrix, not a differential operator, and as such can be applied to vector-valued functions without passing to Fourier space. The downside is that $\mathcal{I}_1^n(s, q, \varepsilon)$ contains the factor $1/\varepsilon$ and, in addition, the unbounded differential operator $\mathcal{B}$, which is disadvantageous for accuracy and stability, respectively. Before we tackle this problem, we show in the following lemma that the terms $\mathcal{I}_2^n(s, q, \varepsilon)$ and $\mathcal{I}_3^n(s, q, \varepsilon)$ are not relevant for the construction of the time integrator.

**Lemma 3.1**
*Under the assumptions of Corollary 2.10, there is a constant $C$ such that*

$$\|\mathcal{I}_2^n(s, q, \varepsilon)\|_W \leq Cs \qquad and \qquad \|\mathcal{I}_3^n(s, q, \varepsilon)\|_W \leq C\varepsilon s$$

*for all $s \in [0, \tau]$.*

*Proof.* The inequalities follow directly from the fact that $\|E(t)\|_W = 1$ for all $t \in \mathbb{R}$ together with (2.25) and (2.27), and in case of $\mathcal{I}_3^n$ from Corollary 2.10. $\qquad \square$

This lemma implies that if we substitute (3.10) into the third term of (3.7) and then omit $\mathcal{I}_2^n(s, q, \varepsilon)$ and $\mathcal{I}_3^n(s, q, \varepsilon)$, this causes only a contribution of $\mathcal{O}(\tau^2(1 + \varepsilon))$ to the local error, because the variable $s$ is integrated from 0 to $\tau$. Omitting the first term $\mathcal{I}_1^n(s, q, \varepsilon)$, however, would formally cause a contribution of $\mathcal{O}(\tau^2/\varepsilon)$ to the local error, which is more than what we can afford.

In order to identify the dominating part of $\mathcal{I}_1^n(s, q, \varepsilon)$, we decompose

$$P^\perp\mathcal{B}q(t_n + r) = P^\perp\mathcal{B}Pq(t_n + r) + P^\perp\mathcal{B}P^\perp q(t_n + r)$$

and note that $P^\perp\mathcal{B}P = P^\perp\mathcal{A}P$ by definition. Since

$$\|\mathcal{B}P^\perp q(t_n + r)\|_W \leq C\varepsilon \qquad \text{for } t_n + r \in [0, t_{\text{end}}]$$

by Corollary 2.10, it follows that

$$\mathcal{I}_1^n(s, q, \varepsilon) = -\frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)P^\perp\mathcal{A}Pq(t_n + r)\,\mathrm{d}r + \mathcal{O}(s). \tag{3.12}$$

The integral still contains evaluations of $q$ at times $t_n + r$ for $r \in [0, s]$. We fix the solution at time $t_n$ such that

$$-\frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)P^\perp\mathcal{A}Pq(t_n + r)\,\mathrm{d}r = -\frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)P^\perp\mathcal{A}Pq(t_n)\,\mathrm{d}r + \mathcal{J}^n(s, q, \varepsilon)$$

with the error term

$$\mathcal{J}^n(s, q, \varepsilon) = -\frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)P^\perp\mathcal{A}\big(Pq(t_n + r) - Pq(t_n)\big)\,\mathrm{d}r.$$

We want to show that the term $\mathcal{J}^n(s, q, \varepsilon)$ is sufficiently small and hence can be neglected in the construction of the integrator. Note that the straightforward bound

$$\left\|\frac{1}{\varepsilon}\int_0^s \exp\left(-\frac{\mathrm{i}(s-r)}{\varepsilon^2}L\right)P^\perp\mathcal{A}\big(Pq(t_n + r) - Pq(t_n)\big)\,\mathrm{d}r\right\|_W$$

$$\leq \frac{s}{\varepsilon}\sup_{r \in [0,s]}\left\|\mathcal{A}\big(Pq(t_n + r) - Pq(t_n)\big)\right\|_W \leq C\frac{s^2}{\varepsilon}$$

following from (3.6) is not sufficient. If $L$ were invertible, then a factor of $\varepsilon^2$ could be gained by integrating by parts, but we know that $L$ has a nontrivial kernel. At this point, the occurrence of the projector $P^\perp$ in the integrand is crucial.

Let $L_\perp$ denote the restriction of $L$ to the subspace $P^\perp \mathbb{C}^m$, i.e.

$$L_\perp : P^\perp \mathbb{C}^m \to P^\perp \mathbb{C}^m, \qquad L_\perp = LP^\perp = \sum_{\ell=m_0+1}^{m} \lambda_\ell(0)\psi_\ell(0)\psi_\ell^*(0) \tag{3.13}$$

with $\lambda_\ell$ and $\psi_\ell$ from Section 2.2. This mapping is invertible with inverse

$$L_\perp^{-1} : P^\perp \mathbb{C}^m \to P^\perp \mathbb{C}^m, \qquad L_\perp^{-1} = \sum_{\ell=m_0+1}^{m} \frac{1}{\lambda_\ell(0)}\psi_\ell(0)\psi_\ell^*(0), \tag{3.14}$$

cf. [5, Section 3.2] and similar in [4, proof of Proposition 3.4]. The definition of $L_\perp$ implies that $L^j P^\perp = L_\perp^j P^\perp$ for every $j \in \mathbb{N}_0$ and hence that

$$E\left(\tfrac{t}{\varepsilon^2}\right) P^\perp = \exp\left(-\frac{\mathrm{i}t}{\varepsilon^2}L_\perp\right) P^\perp \qquad \text{for all } t \in \mathbb{R}. \tag{3.15}$$

This allows us to prove the following bound for $\mathcal{J}^n(s,q,\varepsilon)$.

**Lemma 3.2**
*Let $p_0 \in W^2$. There is a constant $C$ such that*

$$\|\mathcal{J}^n(s,q,\varepsilon)\|_W \leq C\varepsilon s$$

*for all $s \in [0,\tau]$.*

*Proof.* Using (3.15), we can apply integration by parts. This yields

$$\mathcal{J}^n(s,q,\varepsilon) = -\frac{1}{\varepsilon} \int_0^s \exp\left(-\frac{\mathrm{i}(s-r)}{\varepsilon^2}L_\perp\right) P^\perp \mathcal{A}\big(Pq(t_n+r) - Pq(t_n)\big)\,\mathrm{d}r$$

$$= \mathrm{i}\varepsilon L_\perp^{-1} P^\perp \mathcal{A}\big(Pq(t_n+s) - Pq(t_n)\big) - \mathrm{i}\varepsilon \int_0^s L_\perp^{-1} \exp\left(-\frac{\mathrm{i}(s-r)}{\varepsilon^2}L_\perp\right) \partial_t P^\perp \mathcal{A}Pq(t_n+r)\,\mathrm{d}r.$$

Both terms on the right-hand side are in $\mathcal{O}(\varepsilon s)$ due to (3.6) and Lemma 2.11. $\qquad \square$

Combining (3.12) and Lemma 3.2 leads to

$$\mathcal{I}_1^n(s,q,\varepsilon) = -\frac{1}{\varepsilon} \int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) P^\perp \mathcal{A}Pq(t_n)\,\mathrm{d}r + \mathcal{O}(\varepsilon s + s). \tag{3.16}$$

The term of $\mathcal{O}(\varepsilon s + s)$ can be omitted in the construction of the method. The first term on the right-hand side of (3.16) can be computed numerically, but, as mentioned before, it would lead to an instable method because of the differential operator $\mathcal{A}$. To ensure stability, we approximate $\mathcal{A}$ by a filtered version as, e.g., in [10, Section 2.3] and [11, Section 2]. Unfortunately, the default choice

$$\mathcal{A} \approx \frac{\mathrm{i}}{\tau}\sin\left(\frac{\tau}{\mathrm{i}}\mathcal{A}\right) \tag{3.17}$$

used in [10, 11] does *not* work in our case; cf. Remark 3.10 below. As we will see later, the proper choice is

$$\mathcal{A} \approx \frac{\mathrm{i}\varepsilon}{\tau}\sin\left(\frac{\tau}{\mathrm{i}\varepsilon}\mathcal{A}\right) = \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} \tag{3.18}$$

with $\widetilde{\mathcal{A}}_\gamma = \frac{\mathrm{i}}{\gamma}\sin\left(\frac{\gamma}{\mathrm{i}}\mathcal{A}\right).$

**Lemma 3.3**
*For all $\gamma > 0$ and $f \in W$ the inequality*

$$\left\|\widetilde{\mathcal{A}}_\gamma f\right\|_W \leq \frac{1}{\gamma} \|f\|_W \tag{3.19}$$

*holds. Moreover, there is a constant $C$ such that*

$$\left\|\mathcal{A}f - \widetilde{\mathcal{A}}_\gamma f\right\|_W \leq C\gamma \|f\|_{W^2} \tag{3.20}$$

*for all $f \in W^2$ and $\gamma > 0$.*

*Proof.* For every $\gamma > 0$ and every $f \in W$, we obtain by definition of the Wiener algebra that

$$\left\|\widetilde{\mathcal{A}}_\gamma f\right\|_W = \left\|\frac{\mathrm{i}}{\gamma} \sin\left(\frac{\gamma}{\mathrm{i}}\mathcal{A}\right) f\right\|_W = \frac{1}{\gamma}\left\|\sin\left(\gamma A(\cdot)\right) \widehat{f}\right\|_{L^1} \leq \frac{1}{\gamma}\left\|\widehat{f}\right\|_{L^1} = \frac{1}{\gamma}\|f\|_W.$$

If $f \in W^2$, then

$$\left\|\mathcal{A}f - \widetilde{\mathcal{A}}_\gamma f\right\|_W = \left\|\mathrm{i}A(\cdot)\widehat{f} - \frac{\mathrm{i}}{\gamma}\sin\left(\gamma A(\cdot)\right)\widehat{f}\right\|_{L^1} = \int_{\mathbb{R}^d}\left|A(k)\widehat{f}(k) - \frac{1}{\gamma}\sin\left(\gamma A(k)\right)\widehat{f}(k)\right|_2 \mathrm{d}k.$$

For every fixed $k \in \mathbb{R}^d$, Taylor expansion of $s \mapsto \sin\left(\gamma s A(k)\right)$ about $s = 0$ yields

$$\left|A(k)\widehat{f}(k) - \frac{1}{\gamma}\sin\left(\gamma A(k)\right)\widehat{f}(k)\right|_2 \leq \gamma\left|\int_0^1 (1-s)A^2(k)\sin\left(\gamma s A(k)\right)\widehat{f}(k)\,\mathrm{d}s\right|_2 \leq C\gamma|k|_1^2|\widehat{f}(k)|_2,$$

which proves (3.20). □

Next, we estimate the error caused by using $\widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}$ instead of $\mathcal{A}$ in (3.16). As before, the obvious estimate

$$\left\|\frac{1}{\varepsilon}\int_0^s E\left(\frac{s-r}{\varepsilon^2}\right) P^\perp\left(\mathcal{A} - \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}\right)Pq(t_n)\,\mathrm{d}r\right\|_W \leq \frac{s}{\varepsilon}\left\|\left(\mathcal{A} - \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}\right)Pq(t_n)\right\|_W \leq C\frac{s\tau}{\varepsilon^2}\left\|Pq(t_n)\right\|_{W^2}$$

is inadequate for this purpose. Once again, we have to take advantage of the fact that the restriction of $L$ to the subspace $P^\perp \mathbb{C}^m$ is invertible.

**Lemma 3.4**
*There is a constant $C$ such that for every $f \in W^2$*

$$\frac{1}{\varepsilon}\left\|\int_0^s E\left(\frac{s-r}{\varepsilon^2}\right) P^\perp\left(\mathcal{A} - \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}\right)Pf\,\mathrm{d}r\right\|_W \leq C\tau\|f\|_{W^2}.$$

*Proof.* It follows from (3.15) that

$$\int_0^s E\left(\frac{s-r}{\varepsilon^2}\right) P^\perp\,\mathrm{d}r = \int_0^s \exp\left(-\frac{\mathrm{i}(s-r)}{\varepsilon^2}L_\perp\right) P^\perp\,\mathrm{d}r = \frac{\varepsilon^2}{\mathrm{i}}\left(I - \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\right) L_\perp^{-1}P^\perp. \tag{3.21}$$

With (3.20) we obtain

$$\frac{1}{\varepsilon}\left\|\int_0^s E\left(\frac{s-r}{\varepsilon^2}\right) P^\perp\left(\mathcal{A} - \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}\right)Pf\,\mathrm{d}r\right\|_W \leq C\varepsilon\left\|\left(\mathcal{A} - \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}\right)Pf\,\mathrm{d}r\right\|_W \leq C\tau\|f\|_{W^2}.$$

□

By combining (3.10) with Lemmas 3.1, 3.2, and 3.4 we obtain the representation

$$P^\perp q(t_n + s) = E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q(t_n) - \frac{1}{\varepsilon} \int\limits_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) \, \mathrm{d}r \; P^\perp \widetilde{\mathcal{A}}_{\frac{r}{\varepsilon}} P q(t_n) + \mathcal{O}(\tau + \varepsilon\tau).$$

Omitting the term of $\mathcal{O}(\tau + \varepsilon\tau)$ yields an approximation for $P^\perp q(t_n + s)$ which we plug into (3.7). Then, all remaining integrals can be computed analytically. This completes the construction of the desired uniformly accurate time integrator.

In order to write this method in a concise form, we introduce the abbreviation

$$\mathcal{G}_\tau(s) := \frac{1}{\varepsilon} \int\limits_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) \, \mathrm{d}r \; P^\perp \widetilde{\mathcal{A}}_{\frac{r}{\varepsilon}} P. \tag{3.22}$$

With this notation, the numerical method to compute approximations $q^n \approx q(t_n)$ at $t_n = n\tau \in [0, t_{\mathrm{end}}]$ is defined by $q^{n+1} = \Phi_\tau(q^n)$ with numerical flow

$$\Phi_\tau(q^n) := \mathcal{E}_\mathcal{B}\left(\tfrac{\tau}{\varepsilon^2}\right) q^n + \int\limits_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T^{\mathrm{sym}}(Pq^n) \, \mathrm{d}s \tag{3.23}$$

$$+ \frac{1}{2} \int\limits_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T^{\mathrm{sym}}_{\mathrm{perm}}\left(E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(s)q^n, Pq^n, Pq^n\right) \mathrm{d}s.$$

Details concerning implementation and in particular the computation of the integrals will be discussed in Section 4.

We would like to point out that the operators $\mathcal{E}_\mathcal{B}(t)$ and $\mathcal{G}_\tau(s)$ also depend on the parameter $\varepsilon$, but for the sake of clarity, we refrain from expressing this in the notation.

## 3.3 Error analysis

Our next goal is to prove the following error bound, which implies that the time integrator for the REE is uniformly accurate and of order one.

**Theorem 3.5**
*Let $p_0 \in W^2$, let*

$$q \in C^2([0, t_{end}]; W) \cap C^1([0, t_{end}]; W^1) \cap C([0, t_{end}]; W^2)$$

*be the solution of (2.24) and let $q^n = \Phi_\tau^n(p_0)$ be the approximations defined by (3.23) with step size $\tau = t_{end}/N$ for $N \in \mathbb{N}$. Suppose that the numerical solution is uniformly bounded: there is a constant $C$ such that*

$$\sup_{\varepsilon \in (0,1]} \max_{n=0,\dots,N} \|q^n\|_W \leq C \tag{3.24}$$

*for all $\tau$ and $N$. Then, under Assumptions 2.1 and 2.3, the global error is bounded by*

$$\max_{n=0,\dots,N} \|q^n - q(t_n)\|_{L^\infty} \leq C_{global} \, \tau \tag{3.25}$$

*with a constant $C_{global}$ which depends on $t_{end}$, but not on $\varepsilon, N$, and $\tau$.*

**Remark 3.6**
*For sufficiently small $\tau$ uniform boundedness of the numerical solution can be verified with a classical bootstrapping argument as, e.g., in [11, proof of Theorem 3.1].*

The rest of this subsection is devoted to the proof of Theorem 3.5. We start with the following observation.

**Lemma 3.7** (Local error)
*Under the assumptions of Theorem 3.5 there is a constant $C_{local}$ such that*

$$\|q(t_{n+1}) - \Phi_\tau(q(t_n))\|_W \le C_{local}\,\tau^2$$

*for all $n = 0, \ldots, N-1$ and for all $\varepsilon \in (0, 1]$.*

*Proof.* This is an immediate consequence of how we have derived the time integrator in Section 3.2. $\quad\square$

For the proof of stability, we need the following auxiliary result.

**Lemma 3.8**
*For every $f \in W$ there is a constant $C$ such that*

$$\left\| E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp f - \mathcal{G}_\tau(s) f \right\|_W \le 2 \|f\|_W$$

*for $s \in [0, \tau]$.*

*Proof.* We bound the two parts separately. Since the matrix $E\left(\tfrac{s}{\varepsilon^2}\right) = \exp\left(-\tfrac{\mathrm{i}s}{\varepsilon^2}L\right)$ is unitary, it follows that

$$\left\| E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp f \right\|_W \le \left\| P^\perp f \right\|_W \le \|f\|_W \,.$$

For $s \in [0, \tau]$, the second part can be bounded with (3.19) by

$$\|\mathcal{G}_\tau(s) f\|_W = \left\| \frac{1}{\varepsilon} \int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)\,\mathrm{d}r\ P^\perp \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} Pf \right\|_W \le \frac{\tau}{\varepsilon} \left\| \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} Pf \right\|_W \le \|f\|_W \,.$$

$$\square$$

**Proposition 3.9** (Stability)
*Let $f, g \in W$ and $\tau > 0$. Then the numerical method (3.23) satisfies*

$$\|\Phi_\tau(f) - \Phi_\tau(g)\|_W \le \mathrm{e}^{\tau C}\|f - g\|_W. \tag{3.26}$$

*Proof.* From the definition of the numerical flow $\Phi_\tau$ we obtain

$$\Phi_\tau(f) - \Phi_\tau(g)$$

$$= \mathcal{E}_\mathcal{B}\left(\tfrac{\tau}{\varepsilon^2}\right)(f - g) + \int_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \left[T^{\mathrm{sym}}\left(Pf\right) - T^{\mathrm{sym}}\left(Pg\right)\right]\,\mathrm{d}s$$

$$+ \frac{1}{2} \int_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \left[T^{\mathrm{sym}}_{\mathrm{perm}}\left(E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp f - \mathcal{G}_\tau(s)f, Pf, Pf\right) - T^{\mathrm{sym}}_{\mathrm{perm}}\left(E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp g - \mathcal{G}_\tau(s)g, Pg, Pg\right)\right]\,\mathrm{d}s.$$

Since $\mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right)$ is an isometry on $W$ for every $s \in [0, \tau]$, we have

$$\begin{aligned} &\|\Phi_\tau(f) - \Phi_\tau(g)\|_W \\ &\quad\le \|f - g\|_W + \tau \|T^{\mathrm{sym}}(Pf) - T^{\mathrm{sym}}(Pg)\|_W \\ &\quad\quad + \frac{\tau}{2} \sup_{s \in [0,\tau]} \left\| T^{\mathrm{sym}}_{\mathrm{perm}}\left(E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp f - \mathcal{G}_\tau(s)f, Pf, Pf\right) - T^{\mathrm{sym}}_{\mathrm{perm}}\left(E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp g - \mathcal{G}_\tau(s)g, Pg, Pg\right) \right\|_W. \end{aligned}$$

Lemma 3.8 allows us to apply the inequalities (2.26) and (2.28) to infer

$$\|\Phi_\tau(f) - \Phi_\tau(g)\|_W \le \|f - g\|_W + C\tau\,\|f - g\|_W + C\tau \sup_{s \in [0,\tau]} \left\| E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp(f - g) - \mathcal{G}_\tau(s)(f - g) \right\|_W .$$

It follows with Lemma 3.8 that

$$\left\| E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp (f-g) - \mathcal{G}_\tau(s)(f-g) \right\|_W \le C\, \|f-g\|_W$$

for $s \in [0,\tau]$. In conclusion, this yields

$$\|\Phi_\tau(f) - \Phi_\tau(g)\|_W \le \|f-g\|_W + C\tau\, \|f-g\|_W \le \mathrm{e}^{C\tau}\, \|f-g\|_W .$$

<div align="right">□</div>

**Remark 3.10**
*Note that in the proof of Lemma 3.8 and Proposition 3.9 the choice of the filtering is important. If, for instance, we replace $\widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}}$ by the standard choice $\widetilde{\mathcal{A}}_\tau$, then the last inequality in the proof of Lemma 3.8 becomes*

$$\frac{\tau}{\varepsilon} \left\| \widetilde{\mathcal{A}}_\tau P f \right\|_W \le \frac{C}{\varepsilon}\, \|f\|_W .$$

*This is not sufficient for establishing stability in the proof of Proposition 3.9, because with this estimate we would only obtain*

$$\|\Phi_\tau(f) - \Phi_\tau(g)\|_W \le \mathrm{e}^{\tau C/\varepsilon}\|f-g\|_W.$$

*instead of* (3.26).

*Proof of Theorem 3.5.* After these preparations, the proof of Theorem 3.5 follows by combining the local error bound from Lemma 3.7 with the stability result from Proposition 3.9 in the classical construction known as Lady Windermere's fan. <div align="right">□</div>

# 4  Practical implementation and numerical experiments

Recall that numerical approximations are computed by $q^{n+1} = \Phi_\tau(q^n)$. From the representation of the numerical flow $\Phi_\tau$ given in (3.23), however, it is not obvious how this leads to an executable algorithm. In the following subsection, we explain how to compute the integrals in (3.23) analytically. These explanations refer still to the semi-discretization in time without any discretization in space. Then, in Section 4.2, we sketch how to obtain a fully discrete method. Numerical experiments with this method are presented in Section 4.3.

## 4.1  Computation of the integrals

We decompose the numerical flow (3.23) into three parts, i.e. we let $\Phi_\tau = \Phi_{1,\tau} + \Phi_{2,\tau} + \Phi_{3,\tau}$ with

$$\Phi_{1,\tau}(q^n) = \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau}{\varepsilon^2}\right) q^n,$$

$$\Phi_{2,\tau}(q^n) = \int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T^{\mathrm{sym}}(Pq^n)\,\mathrm{d}s,$$

$$\Phi_{3,\tau}(q^n) = \frac{1}{2}\int_0^\tau \mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T^{\mathrm{sym}}_{\mathrm{perm}}\left( E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(s)q^n, Pq^n, Pq^n\right)\,\mathrm{d}s.$$

The abbreviations $\mathcal{E}_{\mathcal{B}}(t)$, $E(t)$, and $\mathcal{G}_\tau(s)$ are repeated here for the convenience of the reader:

$$\mathcal{E}_{\mathcal{B}}(t) = \exp\left(-t(\mathrm{i}L + \varepsilon\mathcal{B})\right),$$

$$E(t) = \exp\left(-\mathrm{i}tL\right),$$

$$\mathcal{G}_\tau(s) = \frac{1}{\varepsilon}\int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right)\,\mathrm{d}r\; P^\perp \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} P.$$

We discuss the three terms $\Phi_{1,\tau}, \Phi_{2,\tau}, \Phi_{3,\tau}$ of the time integrator individually. The exponential $\mathcal{E}_\mathcal{B}\left(\frac{\tau}{\varepsilon^2}\right)$ in the first term $\Phi_{1,\tau}$ is defined via the Fourier transform as in (2.7), but with $\mathcal{B} = \mathcal{A} - c_g \cdot \nabla$ instead of $\mathcal{A}$. Hence, we compute

$$\mathcal{E}_\mathcal{B}\left(\tfrac{\tau}{\varepsilon^2}\right) q^n = \mathcal{F}^{-1}\left( \exp\left( -\frac{\mathrm{i}\tau}{\varepsilon^2}(L + \varepsilon B(\cdot)) \right) \widehat{q}^n(\cdot) \right) =: \mathcal{F}^{-1}\left(\mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau}{\varepsilon^2}\right) \widehat{q}^n\right) \tag{4.1}$$

with $B(k) = A(k) - (c_g \cdot k)I$ and

$$\left(\mathcal{E}_{\mathrm{i}B}\left(t\right) \widehat{q}^n\right)(k) = \exp\left( -\mathrm{i}t(L + \varepsilon B(k)) \right)\widehat{q}^n(k).$$

Note that $B(\mathrm{i}k) = \mathrm{i}B(k)$ since $A(\beta)$ is linear in $\beta$. For the second term $\Phi_{2,\tau}$, we obtain as in (4.1)

$$\mathcal{F}\left( \int_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T^{\mathrm{sym}}(Pq^n) \,\mathrm{d}s \right)(k) = \int_0^\tau \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \,\mathrm{d}s \, \mathcal{F}\left(T^{\mathrm{sym}}(Pq^n)\right)(k).$$

Note that for efficiency reasons it is important to first evaluate the symmetrized nonlinearity in physical space and then take the Fourier transform, as discussed in Section 3.2. The remaining integral of the frequency-dependent matrix exponential $\mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right)$ can be calculated by use of the eigendecomposition

$$L + \varepsilon B(k) = L + B(\varepsilon k) = \Psi(\varepsilon k)\widetilde{\Lambda}(\varepsilon k)\Psi^*(\varepsilon k) \in \mathbb{C}^{m \times m}$$

with the unitary matrix $\Psi(\theta) = \left(\psi_1(\theta)| \ldots |\psi_m(\theta)\right)$, $\theta \in \mathbb{R}^d$, and the diagonal matrix

$$\widetilde{\Lambda}(\theta) = \mathrm{diag}\left(\widetilde{\lambda}_1(\theta), \,\ldots\,, \widetilde{\lambda}_m(\theta)\right), \qquad \widetilde{\lambda}_\ell(\theta) = \lambda_\ell(\theta) - c_g \cdot \theta \quad \text{for } \ell = 1, \ldots, m.$$

Recall that $\lambda_\ell(\theta)$ and $\psi_\ell(\theta)$ are the eigenvalues and eigenvectors, respectively, of $L(\omega, \kappa + \theta) = L + A(\theta)$ as defined in Section 2.2.

Substituting the eigendecomposition into the matrix exponential yields

$$\int_0^\tau \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \,\mathrm{d}s = \int_0^\tau \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2}(L + \varepsilon B(k)) \right) \,\mathrm{d}s = \Psi(\varepsilon k) \int_0^\tau \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2}\widetilde{\Lambda}(\varepsilon k) \right) \,\mathrm{d}s \, \Psi^*(\varepsilon k).$$

In order to represent the integral of the diagonal matrix in a convenient way we introduce the usual function

$$\varphi_1 : \mathbb{C} \to \mathbb{C}, \qquad \varphi_1(z) = \int_0^1 \mathrm{e}^{sz} \,\mathrm{d}s = \int_0^1 \mathrm{e}^{(1-s)z} \,\mathrm{d}s = \begin{cases} \frac{\mathrm{e}^z - 1}{z} & \text{for } z \neq 0, \\ 1 & \text{else.} \end{cases}$$

This results in

$$\int_0^\tau \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2}\widetilde{\Lambda}(\varepsilon k) \right) \,\mathrm{d}s = \tau \,\mathrm{diag}\left( \varphi_1\left( -\tfrac{\mathrm{i}\tau}{\varepsilon^2}\widetilde{\lambda}_1(\varepsilon k) \right), \,\ldots\,, \varphi_1\left( -\tfrac{\mathrm{i}\tau}{\varepsilon^2}\widetilde{\lambda}_m(\varepsilon k) \right) \right) =: M_{\mathrm{int}}^0(\tau, \varepsilon k).$$

Finally, we consider the last term $\Phi_{3,\tau}$ of the numerical flow. Applying the Fourier transform as in the previous two cases gives

$$\mathcal{F}\left( \frac{1}{2} \int_0^\tau \mathcal{E}_\mathcal{B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) T_{\mathrm{perm}}^{\mathrm{sym}}\left( E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(s)q^n, Pq^n, Pq^n \right) \,\mathrm{d}s \right)(k)$$

$$= \frac{1}{2} \int_0^\tau \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \, \mathcal{F}\left(T_{\mathrm{perm}}^{\mathrm{sym}}\left( E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(s)q^n, Pq^n, Pq^n \right) \right)(k) \,\mathrm{d}s. \tag{4.2}$$

In contrast to the nonlinearity in $\Phi_{2,\tau}$, the first argument of $T_{\text{perm}}^{\text{sym}}$ is not independent of the integration variable $s$ and includes a differential operator. It is therefore necessary to proceed in a different way. We start by expressing $\mathcal{G}_\tau(s)$ without an integral. It follows from (3.21) that

$$\mathcal{G}_\tau(s)q^n = \frac{1}{\varepsilon} \int_0^s E\left(\tfrac{s-r}{\varepsilon^2}\right) \mathrm{d}r \, P^\perp \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} P q^n = \frac{\varepsilon}{\mathrm{i}} \left(I - \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\right)(L_\perp)^{-1} P^\perp \widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} P q^n.$$

The filtered differential operator (3.18) applied to $Pq^n$ can be calculated using the Fourier transform:

$$\widetilde{\mathcal{A}}_{\frac{\tau}{\varepsilon}} P q^n = \mathcal{F}^{-1}(\widetilde{A}_{\frac{\tau}{\varepsilon}} P \widehat{q}^n) =: b^n \qquad \text{with} \qquad \widetilde{A}_{\frac{\tau}{\varepsilon}}(k) = \frac{\mathrm{i}\varepsilon}{\tau} \sin\left(\frac{\tau}{\varepsilon}A(k)\right).$$

Substituting the previous calculations in $E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(\tau)q^n$ together with (3.15) yields

$$E\left(\tfrac{s}{\varepsilon^2}\right) P^\perp q^n - \mathcal{G}_\tau(\tau)q^n = \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right) P^\perp q^n + \mathrm{i}\varepsilon\left(I - \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\right)(L_\perp)^{-1} P^\perp b^n$$

$$= \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\left(P^\perp q^n - \mathrm{i}\varepsilon(L_\perp)^{-1}P^\perp b^n\right) + \mathrm{i}\varepsilon(L_\perp)^{-1}P^\perp b^n$$

and hence

$$\mathcal{F}\Big(T_{\text{perm}}^{\text{sym}}\left(E\left(\tfrac{s}{\varepsilon^2}\right)P^\perp q^n - \mathcal{G}_\tau(s)q^n, Pq^n, Pq^n\right)\Big) = \widehat{Y}_1^n(s) + \widehat{Y}_2^n \tag{4.3}$$

with

$$\widehat{Y}_1^n(s) = \mathcal{F}\left(T_{\text{perm}}^{\text{sym}}\left(\exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\left(P^\perp q^n - \mathrm{i}\varepsilon(L_\perp)^{-1}P^\perp b^n\right), Pq^n, Pq^n\right)\right),$$

$$\widehat{Y}_2^n = \mathcal{F}\left(T_{\text{perm}}^{\text{sym}}\left(\mathrm{i}\varepsilon(L_\perp)^{-1}P^\perp b^n, Pq^n, Pq^n\right)\right).$$

After inserting (4.3) in (4.2), we are left with the two integrals

$$\frac{1}{2}\int_0^\tau \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) \widehat{Y}_1^n(s)\,\mathrm{d}s \quad \text{and} \quad \frac{1}{2}\int_0^\tau \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right)\,\mathrm{d}s\,\widehat{Y}_2^n.$$

Since $\widehat{Y}_2^n$ is already independent of the integration variable $s$, the second integral can be treated in the same way as $\Phi_{2,\tau}$.

In order to compute the first integral analytically, we have to cast $\widehat{Y}_1^n(s)$ into a form where the integration variable $s$ does not appear in the first argument of $T_{\text{perm}}^{\text{sym}}$ any more. For this purpose, we use the eigendecomposition

$$L_\perp = \Psi(0)\Lambda(0)\Psi^*(0)$$

with unitary matrix $\Psi(0) = \big(\psi_1(0)|\dots|\psi_m(0)\big)$ and diagonal matrix

$$\Lambda(0) = \operatorname{diag}\Big(0,\ \dots\ ,0,\lambda_{m_0+1}(0),\ \dots\ ,\lambda_m(0)\Big).$$

If the argument of $\Psi$ or $\Lambda$ is equal to zero, this argument will be omitted in subsequent calculations. With the eigendecomposition

$$L_\perp^{-1} = \Psi\Lambda_\perp^{-1}\Psi^* \qquad \text{with} \qquad \Lambda_\perp^{-1} = \operatorname{diag}(0,\dots,0,\lambda_{m_0+1}^{-1},\dots,\lambda_m^{-1})$$

obtained from (3.14), the first entry of the nonlinearity in $\widehat{Y}_1^n(s)$ can be reformulated as

$$\exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}L_\perp\right)\left(P^\perp q^n - \mathrm{i}\varepsilon L_\perp^{-1}P^\perp b^n\right) = \Psi\exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}\Lambda\right)\left(\Psi^*P^\perp q^n - \mathrm{i}\varepsilon\Lambda_\perp^{-1}\Psi^*P^\perp b^n\right)$$

$$= \sum_{\ell=m_0+1}^m \exp\left(-\frac{\mathrm{i}s}{\varepsilon^2}\lambda_\ell\right) a_\ell^n \psi_\ell \tag{4.4}$$

with $a_\ell^n := \psi_\ell^* \left( P^\perp q^n - \frac{\mathrm{i}\varepsilon}{\lambda_\ell} P^\perp b^n \right)$. This yields

$$\widehat{Y}_1^n(s) = \mathcal{F} \left( T_{\mathrm{perm}}^{\mathrm{sym}} \left( \sum_{\ell=m_0+1}^{m} \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) a_\ell^n \psi_\ell, \ Pq^n, \ Pq^n \right) \right). \tag{4.5}$$

At this point, we would like to take the sum and the scalar factor $\exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right)$ out of the nonlinearity. Since $T^{\mathrm{sym}}$ is only real-trilinear (cf. Section 2.1) and the definition of $T_{\mathrm{perm}}^{\mathrm{sym}}$ in (2.22) is based on $T^{\mathrm{sym}}$, we define yet another nonlinearity, namely

$$T_{\mathrm{perm}}(f_1, f_2, f_3) = T(f_1, f_2, f_3) + T(f_1, f_3, f_2) + T(f_2, f_1, f_3) \tag{4.6}$$
$$+ T(f_2, f_3, f_1) + T(f_3, f_1, f_2) + T(f_3, f_2, f_1)$$

for $f_1, f_2, f_3 \in \mathbb{C}^m$. This is similar to (2.22), but with $T^{\mathrm{sym}}$ replaced by $T$. By definition, $T_{\mathrm{perm}} : \mathbb{C}^m \times \mathbb{C}^m \times \mathbb{C}^m \to \mathbb{C}^m$ inherits the property of being trilinear from $T$. With (4.6) we can now write (2.22) in the equivalent form

$$T_{\mathrm{perm}}^{\mathrm{sym}}(f_1, f_2, f_3) = T_{\mathrm{perm}}(f_1, f_2, \overline{f_3}) + T_{\mathrm{perm}}(f_1, \overline{f_2}, f_3) + T_{\mathrm{perm}}(\overline{f_1}, f_2, f_3).$$

Applying this to (4.5) and using trilinearity of $T_{\mathrm{perm}}$ leads to

$$\widehat{Y}_1^n(s) = 2 \sum_{\ell=m_0+1}^{m} \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \mathcal{F} \left( T_{\mathrm{perm}} \left( a_\ell^n \psi_\ell, \ Pq^n, \ \overline{Pq^n} \right) \right)$$
$$+ \sum_{\ell=m_0+1}^{m} \exp\left( +\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \mathcal{F} \left( T_{\mathrm{perm}} \left( \overline{a_\ell^n \psi_\ell}, \ Pq^n, \ Pq^n \right) \right).$$

Substituting this representation into

$$\frac{1}{2} \int_0^\tau \mathcal{E}_{\mathrm{i}B} \left( \tfrac{\tau-s}{\varepsilon^2} \right) \widehat{Y}_1^n(s) \, \mathrm{d}s$$

provides the desired form of the integral, which allows for its analytical calculation. The two parts of $\widehat{Y}_1^n(s)$ have a very similar structure and can thus be treated in essentially the same way. Hence, we will perform the calculation of the integral exemplarily for the first part of $\widehat{Y}_1^n(s)$, i.e.

$$\sum_{\ell=m_0+1}^{m} \int_0^\tau \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \mathcal{E}_{\mathrm{i}B} \left( \tfrac{\tau-s}{\varepsilon^2} \right) \, \mathrm{d}s \, \mathcal{F} \left( T_{\mathrm{perm}} \left( a_\ell^n \psi_\ell, Pq^n, \overline{Pq^n} \right) \right).$$

Inserting again the eigendecomposition $L + B(\varepsilon k) = \Psi(\varepsilon k) \widetilde{\Lambda}(\varepsilon k) \Psi^*(\varepsilon k)$ into the matrix exponential

$$\mathcal{E}_{\mathrm{i}B} \left( \tfrac{\tau-s}{\varepsilon^2} \right) = \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2} (L + B(\varepsilon k)) \right)$$

yields for $\ell = m_0 + 1, \ldots, m$

$$\int_0^\tau \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \mathcal{E}_{\mathrm{i}B} \left( \tfrac{\tau-s}{\varepsilon^2} \right) \, \mathrm{d}s = \Psi(\varepsilon k) \int_0^\tau \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2} \widetilde{\Lambda}(\varepsilon k) \right) \, \mathrm{d}s \, \Psi^*(\varepsilon k).$$

The remaining integral can be expressed by the function $\varphi_1$:

$$\int_0^\tau \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \exp\left( -\frac{\mathrm{i}(\tau-s)}{\varepsilon^2} \widetilde{\Lambda}(\varepsilon k) \right) \, \mathrm{d}s$$

$$= \exp\left( -\frac{\mathrm{i}\tau}{\varepsilon^2} \widetilde{\Lambda}(\varepsilon k) \right) \int_0^\tau \exp\left( -\frac{\mathrm{i}s}{\varepsilon^2} \lambda_\ell \right) \exp\left( \frac{\mathrm{i}s}{\varepsilon^2} \widetilde{\Lambda}(\varepsilon k) \right) \, \mathrm{d}s$$

$$= \tau \operatorname{diag} \left( \exp\left( -\frac{\mathrm{i}\tau}{\varepsilon^2} \widetilde{\lambda}_j(\varepsilon k) \right)_j \varphi_1 \left( \frac{\mathrm{i}\tau}{\varepsilon^2} \left( -\lambda_\ell + \widetilde{\lambda}_j(\varepsilon k) \right) \right)_j \right) =: M_{\mathrm{int},\ell}^-(\tau, \varepsilon k).$$

For the second part of $\widehat{Y}_1^n(s)$, where the sign in the scalar exponential factor differs, we have to replace this by

$$\int_0^\tau \exp\left(\frac{\mathrm{i}s}{\varepsilon^2}\lambda_\ell\right)\exp\left(-\frac{\mathrm{i}(\tau-s)}{\varepsilon^2}\widetilde{\Lambda}(\varepsilon k)\right)\,\mathrm{d}s$$

$$= \tau\,\mathrm{diag}\left(\exp\left(-\frac{\mathrm{i}\tau}{\varepsilon^2}\widetilde{\lambda}_j(\varepsilon k)\right)_j\varphi_1\left(\frac{\mathrm{i}\tau}{\varepsilon^2}\big(\lambda_\ell+\widetilde{\lambda}_j(\varepsilon k)\big)\right)_j\right)=:M_{\mathrm{int},\ell}^+(\tau,\varepsilon k).$$

In summary, the new time integrator can be written without integrals as follows:

$$q^{n+1}(k)=\mathcal{F}^{-1}\widehat{q}^{n+1}(k),$$

$$\widehat{q}^{n+1}(k)=\exp\left(-\frac{\mathrm{i}\tau}{\varepsilon^2}(L+\varepsilon B(k))\right)\widehat{q}^n(k)$$

$$+\Psi(\varepsilon k)M_{\mathrm{int}}^0(\tau,\varepsilon k)\Psi^*(\varepsilon k)\left[\mathcal{F}\left(T^{\mathrm{sym}}(Pq^n)\right)(k)+\frac{1}{2}\widehat{Y}_2^n(k)\right]$$

$$+\sum_{\ell=m_0+1}^m\Psi(\varepsilon k)M_{\mathrm{int},\ell}^-(\tau,\varepsilon k)\Psi^*(\varepsilon k)\,\mathcal{F}\left(T_{\mathrm{perm}}\left(a_\ell^n\psi_\ell,\ Pq^n,\ \overline{Pq^n}\right)\right)(k)$$

$$+\frac{1}{2}\sum_{\ell=m_0+1}^m\Psi(\varepsilon k)M_{\mathrm{int},\ell}^+(\tau,\varepsilon k)\Psi^*(\varepsilon k)\,\mathcal{F}\left(T_{\mathrm{perm}}\left(\overline{a_\ell^n\psi_\ell},\ Pq^n,\ Pq^n\right)\right)(k).$$

Note that all matrices involved depend only on $k$, but not on $n$. In an efficient implementation, these matrices have to be computed only once for all considered values of $k$.

**Remark 4.1**
*In principle, we could also construct a similar time integrator for the envelope equation (2.5a) with the full nonlinearity $T^{sym}(p)=T^{sym}(p,p,p)$ after changing variables as in Section 3.1. The problem is that here the oscillatory part $P^\perp p$ appears in all three arguments of the nonlinearity, and using the decomposition (4.4) with $q$ replaced by $p$ in all three arguments leads to double and triple sums instead of the single sum $\sum_{\ell=m_0+1}^m(\ldots)$. This increases the numerical work per time step significantly and makes the implementation more complicated. On the other hand, Theorems 2.5 and 2.9 show that replacing the envelope equation by the REE in the SVEA does not deteriorate the accuracy. For these reasons, it is preferable to base the construction of the time integrator on the REE.*

## 4.2 Space discretization

In the co-moving coordinate system introduced in Section 3.1, the essential support of the solution remains in a compact subset of $\mathbb{R}^d$. This allows us to consider the REE on a $d$-dimensional cube $[-x_{\max}, x_{\max}]^d$ with periodic boundary conditions instead of the full space $\mathbb{R}^d$. The error caused by this truncation is negligible on finite time intervals as long as $x_{\max}$ is sufficiently large.

Up to now, we have focused on the semi-discretization of the REE in time. In order to compute numerical approximations, however, it is necessary to discretize both space and time. We choose a number $k_{\max} \in \mathbb{N}$ and define the set of multi-indices

$$\mathbb{K}^d=\Big\{k=(k_1,\ \ldots\ ,k_d)\in\mathbb{Z}^d\text{ with }k_j\in\{-k_{\max},\ \ldots\ ,k_{\max}-1\}\Big\}.$$

The method of choice for the space discretization is spectral collocation on the tensor grid

$$h\mathbb{K}^d=\{hk\in\mathbb{T}^d\text{ with }k\in\mathbb{K}^d\}\qquad\text{with mesh size }h=\frac{x_{\max}}{k_{\max}}.$$

This means that the exact solution $q(t_n,\cdot)$ is approximated by a trigonometric polynomial which is either represented by its coefficients or by the function values at the nodes $hk \in h\mathbb{K}^d$. The celebrated fast Fourier transform allows switching between both representations with low computational costs. The

latter representation is convenient for the evaluation of nonlinear terms, whereas the former is used to evaluate functions of differential operators in an efficient way. If, for example, $f$ is the trigonometric polynomial

$$f(x) = \sum_{k \in \mathbb{K}^d} c_k \exp\left(\frac{\mathrm{i}\pi}{x_{\max}} k \cdot x\right)$$

with coefficients $c_k \in \mathbb{R}^d$, then it follows that

$$\mathcal{E}_{\mathcal{B}}\left(\tfrac{\tau-s}{\varepsilon^2}\right) f(x) = \sum_{k \in \mathbb{K}^d} \mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right) c_k \exp\left(\frac{\mathrm{i}\pi}{x_{\max}} k \cdot x\right),$$

where $\mathcal{E}_{\mathrm{i}B}\left(\tfrac{\tau-s}{\varepsilon^2}\right)$ is defined in (4.1). To evaluate the right-hand side, we simply have to apply a matrix exponential to each of the coefficient vectors $c_k$.

## 4.3   Model problem and numerical example

A prominent example of the type (1.1) is the Maxwell–Lorentz system

$$
\begin{aligned}
\partial_t \mathbf{B} &= -\operatorname{curl} \mathbf{E}, \\
\partial_t \mathbf{E} &= \operatorname{curl} \mathbf{B} - \frac{1}{\varepsilon}\mathbf{Q}, \\
\partial_t \mathbf{Q} &= \frac{1}{\varepsilon}(\mathbf{E} - \mathbf{P}) + \varepsilon|\mathbf{P}|_2^2\mathbf{P}, \\
\partial_t \mathbf{P} &= \frac{1}{\varepsilon}\mathbf{Q},
\end{aligned}
\tag{4.7}
$$

cf. [14, 15, 17, 28, 29]. In this system, the Maxwell equations for the electric field $\mathbf{E}(t,x) \in \mathbb{R}^3$ and the magnetic field $\mathbf{B}(t,x) \in \mathbb{R}^3$ are coupled to ordinary differential equations for the polarization $\mathbf{P}(t,x) \in \mathbb{R}^3$ and its time derivative $\frac{1}{\varepsilon}\mathbf{Q}(t,x) \in \mathbb{R}^3$. An insightful derivation of this model and an interpretation of the parameter $\varepsilon$ was given in [17].

In order to illustrate the behavior of our new time integrator, we consider a one-dimensional reduction of the Maxwell–Lorentz system, for which a reference solution can be computed. For the simplification we assume that $\mathbf{E} = (\mathbf{E}_1, \mathbf{E}_2, \mathbf{E}_3)$, $\mathbf{B} = (\mathbf{B}_1, \mathbf{B}_2, \mathbf{B}_3)$, and $\mathbf{P} = (\mathbf{P}_1, \mathbf{P}_2, \mathbf{P}_3)$ depend only on $t$ and $x_1$, but are constant with respect to $x_2$ and $x_3$. Moreover, we assume that

$$\mathbf{B}_1(t,x) = \mathbf{B}_3(t,x) = \mathbf{E}_1(t,x) = \mathbf{E}_2(t,x) = \mathbf{P}_1(t,x) = \mathbf{P}_2(t,x) = \mathbf{Q}_1(t,x) = \mathbf{Q}_2(t,x) = 0 \tag{4.8}$$

for $t = 0$, which implies that (4.8) is true for all $t \geq 0$. These simplifications lead to the model problem

$$
\begin{aligned}
\partial_t \mathbf{B}_2 &= \partial_{x_1} \mathbf{E}_3, && x_1 \in \mathbb{R}, t \in [0, t_{\mathrm{end}}/\varepsilon], \\
\partial_t \mathbf{E}_3 &= \partial_{x_1} \mathbf{B}_2 - \frac{1}{\varepsilon}\mathbf{Q}_3, \\
\partial_t \mathbf{Q}_3 &= \frac{1}{\varepsilon}(\mathbf{E}_3 - \mathbf{P}_3) + \varepsilon|\mathbf{P}_3|^2\mathbf{P}_3, \\
\partial_t \mathbf{P}_3 &= \frac{1}{\varepsilon}\mathbf{Q}_3.
\end{aligned}
$$

For the sake of concise notation, we will henceforth write $x$ instead of $x_1$. If we define the function $u : [0, t_{\mathrm{end}}/\varepsilon] \times \mathbb{R} \to \mathbb{R}^4$ by

$$u(t,x) = \begin{pmatrix} \mathbf{B}_2(t,x) \\ \mathbf{E}_3(t,x) \\ \mathbf{Q}_3(t,x) \\ \mathbf{P}_3(t,x) \end{pmatrix},$$

then the model problem can be interpreted a special case of the nonlinear Friedrichs system (1.1a) with the matrices

$$A_1 = \begin{pmatrix} 0 & -1 & 0 & 0 \\ -1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \end{pmatrix} \quad \text{and} \quad E = \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & -1 & 0 & 1 \\ 0 & 0 & -1 & 0 \end{pmatrix}$$

and the nonlinearity

$$T(f_1, f_2, f_3) = f_1^{\mathrm{T}} \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix} f_2 \begin{pmatrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{pmatrix} f_3.$$

It can be checked that Assumptions 2.3 is fulfilled. The eigenvalues of $L(0, \beta)$ are

$$\omega_1(\beta) = \frac{1}{\sqrt{2}}\sqrt{2 + \beta^2 + \sqrt{4 + \beta^4}}, \qquad \omega_2(\beta) = -\frac{1}{\sqrt{2}}\sqrt{2 + \beta^2 + \sqrt{4 + \beta^4}},$$

$$\omega_3(\beta) = \frac{1}{\sqrt{2}}\sqrt{2 + \beta^2 - \sqrt{4 + \beta^4}}, \qquad \omega_4(\beta) = -\frac{1}{\sqrt{2}}\sqrt{2 + \beta^2 - \sqrt{4 + \beta^4}}$$

for $\beta \in \mathbb{R}$. Since they have constant multiplicities in a neighborhood of $\beta$, one can show that a smooth eigendecomposition exists, cf. [31, Theorem 3.I.1]. In fact, the eigenvalues are even differentiable in zero in this particular example. For Lipschitz continuity, one can verify the boundedness of the derivative with respect to $\beta$.

Throughout, we fix the wave vector $\kappa = 1.2$ and choose $\omega$ to be the largest eigenvalue $\omega_1(\kappa)$ of $L(0, \kappa)$. The initial data are $p_0(x) = e^{-x^2}\psi_1(0)$, where $\psi_1(0)$ is the normalized eigenvector of $L(\omega, \kappa)$ to the eigenvalue $\lambda_1 = 0$. Hence, $p_0$ fulfills the polarization condition (Assumption 2.1).

For our numerical computations we use the change of variables from Section 3.1 with $t_{\mathrm{end}} = 1$. The space discretization is done as explained in Section 4.2 with $x_{\max} = 4$ and $k_{\max} = 32$. Reference solutions for the rescaled versions of the envelope equation and the REE are computed by Strang splitting with an appropriately small step size $\tau_{\mathrm{ref}} \approx 7 \cdot 10^{-7}$; cf. [25, Chapter IV], [7, Chapter 6], and [8] for overviews on splitting methods for PDEs. In both equations, the linear part reads

$$\partial_t f = -\left(\frac{\mathrm{i}}{\varepsilon^2}L - \frac{1}{\varepsilon}\mathcal{B}\right) f$$

and can be propagated exactly in Fourier space; cf. Section 4.1. To approximate the evolution of the nonlinear part, we use Heun's method; cf. [21, Section 5.6.3].

In our first experiment, we test the accuracy of our new time integrator for the REE (3.2) for three different values of $\varepsilon$. In order to confirm the error bound stated in Theorem 3.5, we compute the discrete counterpart of the left-hand side of (3.25), which means that the semi-discrete approximation $q^n$ and the exact solution $q(t_n)$ are replaced by the fully discrete approximation and by the reference solution, respectively. Since these objects are computed on finitely many grid points, the norm $\|\cdot\|_{L^\infty}$ is replaced by the discrete maximum norm. The result is visualized by the blue stars in Figure 1. In each panel of Figure 1, we also depict the corresponding error of the Strang splitting (green crosses). The dashed black lines are reference lines for convergence order one. We find that the error of the Strang splitting is completely erratic for $\tau > C\varepsilon^2$, and that convergence only kicks in for very small step sizes $\tau \leq C\varepsilon^2$. This behavior is to be expected in view of the classical convergence theory for splitting methods, and it demonstrates that splitting and other traditional methods are not suitable for the REE when $\varepsilon$ is small. In contrast, the plots show that our new time integrator converges with order one without any such step size restriction, and that the size of the error is independent of the parameter $\varepsilon$, in accordance with Theorem 3.5. However, the blue stars do not really lie on a straight line, which indicates that actually the observed convergence order fluctuates a bit around the value 1. This is why two reference lines for
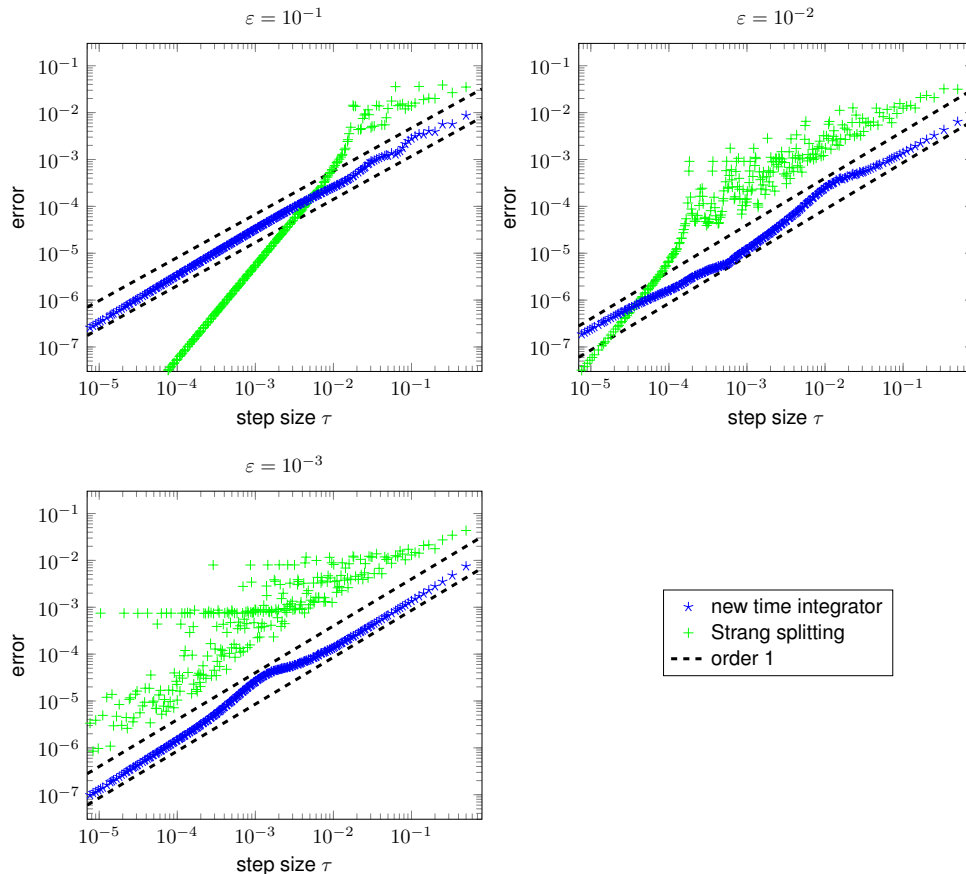
Figure 1: Accuracy of the new time integrator (blue stars) for $\varepsilon = 0.1$ (top left), $\varepsilon = 0.01$ (top right) and $\varepsilon = 0.001$ (bottom left). In addition, the error for Strang splitting (green crosses) is shown for the same values of $\varepsilon$. The dashed black lines are reference lines for order one.

order 1 are depicted in all panels of Figure 1. We conjecture that this behavior is caused by different error terms dominating in different regimes, and we emphasize that the numerical results are clearly not in contradiction to the error bound presented in Theorem 3.5.

Finally, we check the convergence of the REE to the full envelope equation for $\varepsilon \to 0$, as predicted by Theorem 2.9. For simplicity, we use the norm $\| \cdot \|_{L^\infty}$ instead of $\| \cdot \|_{W^1}$ in (2.30), and as before, the exact solutions $p$ and $q$ are replaced by suitable fully discrete reference solutions, which are computed for seven different values of $\varepsilon$. The result is displayed in Figure 2. Comparing the error with the black dashed reference line confirms that replacing the envelope equation by the REE indeed causes an error proportional to $\varepsilon^2$.

# References

[1] W. Bao and C. Su. Uniform error bounds of a finite difference method for the Klein-Gordon-Zakharov system in the subsonic limit regime. *Math. Comp.*, 87(313):2133–2158, 2018. `doi:10.1090/mcom/3278`.

[2] W. Bao and X. Zhao. A uniformly accurate (UA) multiscale time integrator Fourier pseudospectral method for the Klein-Gordon-Schrödinger equations in the nonrelativistic limit regime. *Numer. Math.*, 135(3):833–873, 2017. `doi:10.1007/s00211-016-0818-x`.
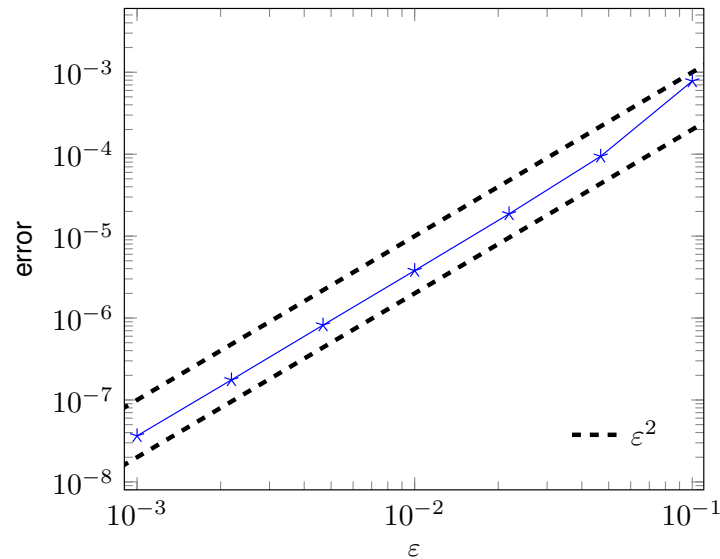
Figure 2: Difference between solution $p$ of the envelope equation and the solution $q$ of the REE for different values of $\varepsilon$.

[3] J. Baumstark. *High-frequency wave-propagation: error analysis for analytical and numerical approximations*. PhD thesis, Karlsruhe Institute of Technology, 2022. `doi:10.5445/IR/1000149719`.

[4] J. Baumstark and T. Jahnke. Improved error bounds for approximations of high-frequency wave propagation in nonlinear dispersive media. CRC 1173 Preprint 2024/5, Karlsruhe Institute of Technology, February 2024. URL: `https://www.waves.kit.edu/downloads/CRC1173_Preprint_2024-5.pdf`, `doi:10.5445/IR/1000168785`.

[5] J. Baumstark, T. Jahnke, and C. Lubich. Polarized high-frequency wave propagation beyond the nonlinear Schrödinger approximation. *SIAM J. Math. Anal.*, 56(1):454–473, 2024. `doi:10.1137/22M1504810`.

[6] S. Baumstark, E. Faou, and K. Schratz. Uniformly accurate exponential-type integrators for Klein-Gordon equations with asymptotic convergence to the classical NLS splitting. *Math. Comp.*, 87(311):1227–1254, 2018. `doi:10.1090/mcom/3263`.

[7] S. Blanes and F. Casas. *A concise introduction to geometric numerical integration*. Monographs and Research Notes in Mathematics. CRC Press, Boca Raton, FL, 2016.

[8] S. Blanes, F. Casas, and A. Murua. Splitting methods for differential equations. *Acta Numer.*, 33:1–161, 2024. `doi:10.1017/S0962492923000077`.

[9] M. Cabrera Calvo and K. Schratz. Uniformly accurate low regularity integrators for the Klein-Gordon equation from the classical to nonrelativistic limit regime. *SIAM J. Numer. Anal.*, 60(2):888–912, 2022. `doi:10.1137/21M1415030`.

[10] Y. Cai and Y. Wang. Uniformly accurate nested Picard iterative integrators for the nonlinear Dirac equation in the nonrelativistic regime. *Multiscale Model. Simul.*, 20(1):164–187, 2022. `doi:10.1137/20M133573X`.

[11] Y. Cai and X. Zhou. Uniformly accurate nested Picard iterative integrators for the Klein-Gordon equation in the nonrelativistic regime. *J. Sci. Comput.*, 92(2):Paper No. 53, 28, 2022. `doi:10.1007/s10915-022-01909-5`.

[12] F. Castella, P. Chartier, F. Méhats, and A. Murua. Stroboscopic averaging for the nonlinear Schrödinger equation. *Found. Comput. Math.*, 15(2):519–559, 2015. `doi:10.1007/s10208-014-9235-7`.

[13] P. Chartier, M. Lemou, F. Méhats, and X. Zhao. Derivative-free high-order uniformly accurate schemes for highly oscillatory systems. *IMA J. Numer. Anal.*, 42(2):1623–1644, 2022. `doi:10.1093/imanum/drab014`.

[14] M. Colin and D. Lannes. Short Pulses Approximations in Dispersive Media. *SIAM Journal on Mathematical Analysis*, 41(2):708–732, 2009. `doi:10.1137/070711724`.

[15] T. Colin, G. Gallice, and K. Laurioux. Intermediate models in nonlinear optics. *SIAM Journal on Mathematical Analysis*, 36(5):1664–1688, 2005. `doi:10.1137/S0036141003423065`.

[16] M. Condon, A. Iserles, K. Kropielnicka, and P. Singh. Solving the wave equation with multifrequency oscillations. *J. Comput. Dyn.*, 6(2):239–249, 2019. `doi:10.3934/jcd.2019012`.

[17] P. Donnat and J. Rauch. Modeling the Dispersion of Light. In *Singularities and Oscillations*, pages 17–35, New York, NY, 1997. Springer New York. `doi:10.1007/978-1-4612-1972-9_2`.

[18] E. Faou and K. Schratz. Asymptotic preserving schemes for the Klein-Gordon equation in the non-relativistic limit regime. *Numer. Math.*, 126(3):441–469, 2014. `doi:10.1007/s00211-013-0567-z`.

[19] B. García-Archilla, J. M. Sanz-Serna, and R. D. Skeel. Long-time-step methods for oscillatory differential equations. *SIAM J. Sci. Comput.*, 20(3):930–963, 1999. `doi:10.1137/S1064827596313851`.

[20] L. Gauckler. Error analysis of trigonometric integrators for semilinear wave equations. *SIAM J. Numer. Anal.*, 53(2):1082–1106, 2015. `doi:10.1137/140977217`.

[21] W. Gautschi. *Numerical Analysis*. SpringerLink. Birkhäuser Boston, Boston, second edition, 2012. `doi:10.1007/978-0-8176-8259-0`.

[22] V. Grimm and M. Hochbruck. Error analysis of exponential integrators for oscillatory second-order differential equations. *J. Phys. A*, 39(19):5495–5507, 2006. `doi:10.1088/0305-4470/39/19/S10`.

[23] E. Hairer, C. Lubich, and Y. Shi. Large-stepsize integrators for charged-particle dynamics over multiple time scales. *Numer. Math.*, 151(3):659–691, 2022. `doi:10.1007/s00211-022-01298-9`.

[24] M. Hochbruck and C. Lubich. A Gautschi-type method for oscillatory second-order differential equations. *Numer. Math.*, 83(3):403–426, 1999. `doi:10.1007/s002110050456`.

[25] W. Hundsdorfer and J. Verwer. *Numerical solution of time-dependent advection-diffusion-reaction equations*, volume 33 of *Springer Series in Computational Mathematics*. Springer-Verlag, Berlin, 2003. `doi:10.1007/978-3-662-09017-6`.

[26] T. Jahnke. Long-time-step integrators for almost-adiabatic quantum dynamics. *SIAM J. Sci. Comput.*, 25(6):2145–2164, 2004. `doi:10.1137/S1064827502411316`.

[27] T. Jahnke and M. Mikl. Adiabatic exponential midpoint rule for the dispersion-managed nonlinear Schrödinger equation. *IMA J. Numer. Anal.*, 39(4):1818–1859, 2019. `doi:10.1093/imanum/dry045`.

[28] J. L. Joly, G. Metivier, and J. Rauch. Global Solvability of the Anharmonic Oscillator Model from Nonlinear Optics. *SIAM Journal on Mathematical Analysis*, 27(4):905–913, 1996. `doi:10.1137/S0036141094273672`.

[29] D. Lannes. High-frequency nonlinear optics: from the nonlinear Schrödinger approximation to ultrashort-pulses equations. *Proceedings of the Royal Society of Edinburgh Section A: Mathematics*, 141(2):253–286, 2011. `doi:10.1017/S030821050900002X`.

[30] C. Lubich and Y. Shi. On a large-stepsize integrator for charged-particle dynamics. *BIT*, 63(1):Paper No. 14, 17, 2023. `doi:10.1007/s10543-023-00951-5`.

[31] J. Rauch. *Hyperbolic partial differential equations and geometric optics*, volume 133 of *Graduate studies in mathematics*. American Mathematical Society, Providence, RI, 2012. `doi:10.1090/gsm/133`.

[32] G. Schneider and H. Uecker. *Nonlinear PDEs : a dynamical systems approach.* Graduate studies in mathematics ; 182. American Mathematical Society, Providence, Rhode Island, 2017.

[33] B. Wang and Y. Jiang. Improved uniform error bounds on parareal exponential algorithm for highly oscillatory systems. *BIT*, 64(1):Paper No. 6, 34, 2024. doi:10.1007/s10543-023-01005-6.

[34] B. Wang and X. Zhao. Geometric two-scale integrators for highly oscillatory system: uniform accuracy and near conservations. *SIAM J. Numer. Anal.*, 61(3):1246–1277, 2023. doi:10.1137/21M1462908.