

HIGH-ORDER UNIFORMLY ACCURATE TIME INTEGRATORS FOR SEMILINEAR WAVE EQUATIONS OF KLEIN–GORDON TYPE IN THE NON-RELATIVISTIC LIMIT

HAIDAR MOHAMAD AND MARCEL OLIVER

ABSTRACT. We introduce a family of high-order time semi-discretizations for semilinear wave equations of Klein–Gordon type with arbitrary smooth nonlinearities that are uniformly accurate in the non-relativistic limit where the speed of light goes to infinity. Our schemes do not require pre-computations that are specific to the nonlinearity and have no restrictions in step size. Instead, they rely upon a general oscillatory quadrature rule developed in a previous paper (Mohamad and Oliver, *SIAM J. Num. Anal.* 59, 2021, 2310–2319).

1. INTRODUCTION

Let X be a Hilbert space with inner product $\langle \cdot, \cdot \rangle$ and associated norm $\|\cdot\|$. We study a semilinear wave equation on X ,

$$c^{-2} \partial_{tt}\phi + L\phi + c^2 \phi = f(\phi, t), \quad (1a)$$

$$\phi(0) = \phi_0, \quad (1b)$$

$$\partial_t \phi(0) = \phi'_0, \quad (1c)$$

where $\phi: [0, T] \rightarrow X$, L is a closed, densely defined, self-adjoint, non-negative operator on X with domain $\mathcal{D}(L)$, c is a positive constant, and $f: \mathcal{D}(L) \times [0, T] \rightarrow X$ a smooth function. Such equations arise, for example, in acoustics, electromagnetics, quantum mechanics, and geophysical fluid dynamics, both in the “relativistic” ($c = 1$) and “nonrelativistic” ($c \gg 1$) regimes. The motivation for studying (1) as written is that it covers two well-studied special cases:

- (i) $X = H^r(\mathbb{T}^d)$, $L = \Delta$, and $f(\phi, t) = |\phi|^2 \phi$, which corresponds to the standard semilinear Klein–Gordon equation.
- (ii) $X = \mathbb{R}^{2d}$, $L = 0$, and $f(\phi, t) = -2e^{-c^2 t J} \nabla V(e^{c^2 t J} \phi)$, where J is the canonical symplectic matrix in $2d$ dimensions and V is a smooth potential. Changing variables via $q = e^{c^2 t J} \phi$, we can write the system in the standard form

$$\dot{q} = p, \quad (2a)$$

$$\frac{1}{2c^2} \dot{p} = Jp - \nabla V(q). \quad (2b)$$

This system has been used as a finite-dimensional toy model for rotating fluid flow, where the limit $c \rightarrow \infty$ corresponds to a rapidly rotating earth.

Date: January 7, 2024.

2020 Mathematics Subject Classification. Primary 65L11; Secondary 35Q40, 35B25.

Key words and phrases. Semilinear wave equations, oscillatory problems, high-frequency limit, oscillatory integrator.

Analytically, the non-relativistic limit regime is well-studied for these two examples. We refer the reader to [11, 18] for the case of the Klein–Gordon equation and to [10, 14, 15] for the case of system (2). Numerically, equation (1) is extensively studied in the relativistic regime [12, 21]. However, due to the high oscillatory character of the solutions when c is large, most numerical methods suffer from severe time step restriction in the non-relativistic regime.

Several authors have considered the problem of finding “asymptotics-preserving numerical schemes”, i.e., schemes that perform uniformly in this singular limit. Some of these schemes [11] are based on a modulated Fourier expansion of the exact solution [9, 16] where the highly oscillatory problem in (2) is reduced to a non-oscillatory limit Schrödinger equation for which no time step restriction is needed. Other schemes are based on multiscale expansions of the exact solution [3, 6]. Chartier *et al.* [7] recently introduced a new method which employs an averaging transformation to soften the stiffness of the problem, hence allowing standard schemes to retain their order of convergence. Baumstark *et al.* [4] construct first and second order uniformly accurate integrators for the Klein–Gordon equation with *cubic* nonlinearity by integrating the trigonometric products arising from a suitable mild formulation explicitly.

In this paper, we develop a family of high-order asymptotics-preserving schemes for (1) that do not require pre-computations tied to the specific nonlinearity f and have no restrictions in time step size. The construction of the new schemes is explained in Section 4. We outline here their main ingredients, where the first two follow the prior work [4]:

- (i) Reformulate (1) as a coupled first order system using a linear transformation.
- (ii) Factor out the rapidly rotating phase to make it explicit.
- (iii) Iterate the resulting mild formulation up to the desired order for the coupled first order system in the new variables.
- (iv) Use the quadrature rule developed in [19] to handle the high oscillatory integral in the resulting mild formulation and complete the construction of the scheme.

The remainder of the paper is structured as follows. In Section 2, we state some properties of the operator L within the framework of its functional calculus. In Section 3, we introduce quadrature rules for the approximation of highly oscillatory Banach-space-valued functions in specific settings that will fit the construction of our schemes. Section 4 is devoted to the detailed construction of the schemes. Our main result on the order of convergence of the schemes is stated and proved in Section 5. In Section 6, we demonstrate that the new schemes are accurate to their expected order and that their error behavior is indeed uniform in c .

2. PRELIMINARIES

We define the operators $B_c = c^{-1} \sqrt{L + c^2}$ and $A_c = c^2 B_c - c^2$. These operators are well-defined via the spectral theorem for densely defined normal operators (e.g. [20]). Indeed, for any densely defined normal operator $P: \mathcal{D}(P) \subseteq X \rightarrow X$, there exists a unique spectral measure E_P on the Borel σ -algebra $\mathcal{B}(\mathbb{C})$ into the orthogonal

projections on X such that

$$P = \int_{\mathbb{C}} \lambda \, dE_P(\lambda) = \int_{\sigma(P)} \lambda \, dE_P(\lambda). \quad (3)$$

This integral representation of P allows us to define the assignments $P \mapsto f(P)$ for any E -a.e. finite measurable function f by the formula

$$f(P) = \int_{\sigma(P)} f(\lambda) \, dE_P(\lambda) \quad (4)$$

with domain

$$\mathcal{D}(f(P)) = \left\{ x \in X : \int_{\sigma(P)} |f(\lambda)|^2 \, d\langle E_P(\lambda)x, x \rangle < \infty \right\}. \quad (5)$$

Definition 1. Let A, B be two densely defined normal operators. If $\mathcal{D}(AB) \subseteq \mathcal{D}(BA)$ and $AB = BA$ on $\mathcal{D}(AB)$, we write $AB \subseteq BA$ and say that “ A commutes with B .”

We fix in what follows an operator $J \in \mathcal{L}(X)$ such that

$$JL \subseteq LJ, \quad J^* = -J, \quad \text{and} \quad J^2 = -I. \quad (6)$$

We now collect important elementary properties of the operators $J, A_c,$ and B_c .

Lemma 2. *The operators $J, A_c,$ and B_c satisfy the following properties.*

- (i) $\mathcal{D}(L) \subseteq \mathcal{D}(A_c) = \mathcal{D}(B_c) = \mathcal{D}(JA_c) = \mathcal{D}(JB_c),$
- (ii) $\|A_c u\|_{\mathcal{D}(L^j)} \leq \frac{1}{2} \|u\|_{\mathcal{D}(L^{j+1})}$ for any $j \in \mathbb{N},$
- (iii) J and e^{tJ} commute with $f(L)$ for any measurable function $f: \mathbb{R} \rightarrow \mathbb{R}$ and $t \in \mathbb{R};$ in particular, J and e^{tJ} commute with $A_c, B_c,$ and $B_c^{-1},$
- (iv) e^{tJA_c} commutes with J and $f(L)$ for any measurable function $f: \mathbb{R} \rightarrow \mathbb{R}$ and $t \in \mathbb{R},$
- (v) $\|e^{tJA_c}\| \leq 1,$ and
- (vi) $\|(e^{tJA_c} - I)u\| \leq \frac{1}{2} |t| \|u\|_{\mathcal{D}(L)}.$

Proof. The inclusion in (i) follows directly from

$$\int_{\sigma(L)} |\lambda + c^2| \, d\langle E_L(\lambda)u, u \rangle \leq (c^2 + \frac{1}{2}) \|u\|^2 + \frac{1}{2} \int_{\sigma(L)} |\lambda|^2 \, d\langle E_L(\lambda)u, u \rangle; \quad (7)$$

the remaining identities are obvious. To prove (ii), we note that, for $\lambda \geq 0$

$$c\sqrt{\lambda + c^2} - c^2 \leq \frac{\lambda}{2}, \quad (8)$$

and

$$\|A_c u\|_{\mathcal{D}(L^j)}^2 = \int_{\sigma(L)} \left| c\sqrt{\lambda + c^2} - c^2 \right|^2 (1 + |\lambda|^2)^j \, d\langle E_L(\lambda)u, u \rangle. \quad (9)$$

For (iii), we recall that J is bounded and commutes with L . Thus, by [20, Proposition 5.15], $J E_L(K) = E_L(K) J$ for all $K \in \mathcal{B}(\mathbb{C})$. Consequently, $e^{tJ} E_L(K) = E_L(K) e^{tJ}$ for all $K \in \mathcal{B}(\mathbb{C})$ and $t \in \mathbb{R}$. Then the claim is a direct consequence of [20, Proposition 4.23]. For (iv), note that

$$e^{tJA_c} = \int_{\mathbb{C}^2} e^{ct\lambda(\sqrt{\mu^2 + c^2} - c)} \, dE_J(\lambda) \, dE_L(\mu) \quad (10)$$

where the integral is with respect to the product measure $E_J \otimes E_L(K_1 \times K_2) = E_J(K_1) E_L(K_2)$ for all $K_1, K_2 \in \mathcal{B}(\mathbb{C})$. Hence, E_J and E_L commute with $E_J \otimes E_L$

in the sense that $E_J \otimes E_L(K_1 \times K_2) E(K_3) = E(K_3) E_J \otimes E_L(K_1 \times K_2)$ for all $K_1, K_2, K_3 \in \mathcal{B}(\mathbb{C})$. Once again, the claim follows from [20, Proposition 4.23]. Estimate (v) is a direct consequence of the skew-symmetry of J . Finally, to prove estimate (vi), let $u \in \mathcal{D}(L)$. Since the spectrum of J is purely imaginary and $|e^{ix} - 1|^2 \leq x^2$ for $x \in \mathbb{R}$, we estimate

$$\begin{aligned} \|(e^{tJA_c} - I)u\|^2 &= \int_{\sigma(J) \times \sigma(L)} |e^{tc\lambda(\sqrt{\mu^2 + c^2} - c)} - 1|^2 \langle dE_J(\lambda) dE_L(\mu)u, u \rangle \\ &\leq t^2 \int_{\sigma(J) \times \sigma(L)} |c\lambda(\sqrt{\mu^2 + c^2} - c)|^2 \langle dE_J(\lambda) dE_L(\mu)u, u \rangle \\ &\leq t^2 \|JA_c u\|^2. \end{aligned} \quad (11)$$

The claim then follows by estimate (ii). \square

Remark 1. Lemma 2(iii) and (iv) imply that if P and Q are two operators such that P is bounded and $PQ \subseteq QP$, then $\mathcal{D}(PQ) = \mathcal{D}(Q)$ and $P(\mathcal{D}(Q)) \subseteq \mathcal{D}(Q)$. In other words, the domain of Q is invariant under any *bounded* operator commuting with Q . In this paper, the analysis of the numerical schemes assumes solutions of (1) in $\mathcal{D}(L)$ which is, in view of this remark, invariant under any bounded operator commuting with L , in particular J , e^{tJ} , and e^{tJA_c} .

3. QUADRATURE FOR BANACH-SPACE-VALUED FUNCTIONS

In this section, let $(X, \|\cdot\|)$ be a complex Banach space and $\Omega \subset \mathbb{C}$ be open. A function $F: \Omega \rightarrow X$ is analytic if it is differentiable, i.e., provided for every $z_0 \in \Omega$ there exists $F'(z_0) \in X$ such that

$$F'(z_0) = \lim_{z \rightarrow z_0} \frac{F(z) - F(z_0)}{z - z_0}. \quad (12)$$

The following simple lemma shows that estimates on the quadrature error for differentiable complex-valued functions directly imply a corresponding estimate for X -valued functions.

Lemma 3. *Let I be an open interval on the real line and μ a measure on I , possibly discrete. Suppose a quadrature rule with nodes $x_k \in I$ and weights ω_k , $k = 1, \dots, n$ satisfies the error estimate*

$$\left| \int_I f(x) d\mu(x) - \sum_{k=1}^n \omega_k f(x_k) \right| \leq C(n, I) \sup_{x \in I} |f^{(p)}(x)|, \quad (13)$$

for some $p \in \mathbb{N}$ and every $f \in C^p(I, \mathbb{C})$. Then the quadrature rule satisfies the error estimate

$$\left\| \int_I F(x) d\mu(x) - \sum_{k=1}^n \omega_k F(x_k) \right\| \leq C(n, I) \sup_{x \in I} \|F^{(p)}(x)\|, \quad (14)$$

where the integral is understood in the Bochner-sense, for every $F \in C^p(I, X)$.

Proof. Fix $\psi \in X^*$. Let

$$e_n = \int_I F(x) d\mu(x) - \sum_{k=1}^n \omega_k F(x_k). \quad (15)$$

Due to the properties of the Bochner integral,

$$\psi(e_n) = \int_a^b \psi \circ F(x) \, d\mu - \sum_{k=1}^n \omega_k \psi \circ F(x_k), \quad (16)$$

so that, applying (13) to $f = \psi \circ F$, we obtain

$$\begin{aligned} |\psi(e_n)| &\leq C(n, I) \sup_{x \in I} \left| \frac{d^p}{dx^p} [\psi \circ F](x) \right| \\ &\leq C(n, I) \|\psi\|_* \sup_{x \in I} \|F^{(p)}(x)\|. \end{aligned} \quad (17)$$

By the Hahn–Banach theorem, we can choose $\psi \in X^*$ with $\|\psi\|_* \leq 1$ such that $\psi(e_n) = \|e_n\|$. This implies (14). \square

With the help of this lemma, we lift three known estimates for the quadrature error of complex-valued functions to the Banach space setting. The first concerns the trapezoidal rule approximation for the integral of a 1-periodic function F , namely the uniformly weighted Riemann sum

$$T_n(F) = \frac{1}{n} \sum_{k=0}^{n-1} F\left(\frac{k}{n}\right). \quad (18)$$

For given $a > 0$, let

$$\Omega_a = \{z \in \mathbb{C}: -a < \operatorname{Im} z < a\}. \quad (19)$$

Then the following estimate, proved for $X = \mathbb{C}$ in [22], holds true.

Theorem 4. *Let F be an X -valued function, 1-periodic on the real line, analytic with $\|F(z)\| \leq A$ on the strip Ω_a for some $a > 0$. Then for any $n \in \mathbb{N}$,*

$$\left\| \int_0^1 F(x) \, dx - T_n(F) \right\| \leq \frac{2A}{e^{an} - 1}. \quad (20)$$

The constant 2 is as small as possible.

The second concerns the Gauss formula for the integral of a function f defined on the interval $[-1, 1]$,

$$G_m(f) = \sum_{k=1}^m \omega_k f(\xi_k), \quad (21)$$

where the ξ_k are the zeros of the Legendre polynomial p_m of degree m and the weights are given by

$$\omega_k = \frac{2}{(1 - \xi_k^2) [p'_m(\xi_k)]^2}. \quad (22)$$

For given $b > a$ and $\rho > \frac{1}{2}(b - a)$, let $E_\rho(a, b)$ denote the ellipse with foci a, b such that the lengths of its minor and major semi-axes sum up to ρ . Namely,

$$E_\rho(a, b) = \left\{ z \in \mathbb{C}: z = \frac{1}{2}(\rho e^{i\theta} + \frac{1}{4}(b - a)^2 \rho^{-1} e^{-i\theta}) + \frac{1}{2}(a + b), 0 \leq \theta < 2\pi \right\}, \quad (23)$$

and $\Sigma_\rho(a, b)$ the open region in \mathbb{C} bounded by $E_\rho(a, b)$.

The formula (21) can easily be written out for functions defined on an arbitrary interval $[a, b]$ using the affine change of variables

$$\ell: \Sigma_{\frac{2\rho}{b-a}}(-1, 1) \rightarrow \Sigma_\rho(a, b), \quad \ell(x) = \frac{b-a}{2}(x+1) + a. \quad (24)$$

Theorem 5. Fix $k \in \mathbb{N}$, $\varepsilon_0 > 0$, and $\rho > \frac{1}{2}(b-a)$. Set $\alpha = \varepsilon_0 \min\{0, a, b\}$ and $\beta = \varepsilon_0 \max\{0, a, b\}$. Let $F: [\alpha, \beta] \times \Sigma_\rho(a, b) \rightarrow X$ be such that $\zeta \mapsto F(\zeta, z)$ is k -times differentiable for any $z \in \Sigma_\rho(a, b)$ and that $z \mapsto \partial_1^i F(\zeta, z)$, where ∂_1 denotes the partial derivative with respect to the first argument, is analytic on $\Sigma_\rho(a, b)$ for $i = 0, \dots, k-1$ and any $\zeta \in [\alpha, \beta]$ with

$$\max_i \sup_{[\alpha, \beta] \times \Sigma_\rho(a, b)} \|\partial_1^i F\| \leq A_{\text{an}}, \quad (25a)$$

$$\sup_{[\alpha, \beta] \times [a, b]} \|\partial_1^k F\| \leq A_{\text{dif}}. \quad (25b)$$

We abbreviate $f(x) = F(\varepsilon x, x)$. Then, for any $m \in \mathbb{N}$ and $\varepsilon \in (0, \varepsilon_0]$,

$$\left\| \int_a^b f(x) dx - G_m(f) \right\| \leq \frac{16 A_{\text{an}} e^{\varepsilon \rho} \rho^2}{(2\rho - b + a)} \left(\frac{b-a}{2\rho} \right)^{2m+1} + \frac{2 A_{\text{dif}} (b-a)^{k+1} \varepsilon^k}{k!}, \quad (26)$$

where

$$G_m(f) = \frac{b-a}{2} \sum_{i=1}^m \omega_i f(\eta_i) \quad (27)$$

with nodes $\eta_i = \ell(\xi_i)$.

Proof. Writing the Taylor series with respect to the first variable of F , we find that for every $x \in [a, b]$ there exists $\xi = \xi(\varepsilon, x) \in [a, b]$ such that

$$f(x) = \sum_{i=0}^{k-1} \frac{(x-a)^i \varepsilon^i}{i!} \partial_1^i F(\varepsilon a, x) + \frac{(x-a)^k \varepsilon^k}{k!} \partial_1^k F(\varepsilon \xi, x). \quad (28)$$

Thus, the following estimate, proved for $X = \mathbb{C}$ in [8], holds true for the quadrature formula (21) applied on each $f_i(z) = (z-a)^i \partial_1^i F(\varepsilon a, z)$, $i = 0, \dots, k-1$, which is analytic and bounded on $\Sigma_\rho(a, b)$:

$$\begin{aligned} & \left\| \int_a^b f_i(x) dx - G_m(f_i) \right\| \\ & \leq \frac{16 \rho^2}{(2\rho - b + a)} \left(\frac{b-a}{2\rho} \right)^{2m+1} \sup_{z \in \Sigma_\rho(a, b)} \|(z-a)^i \partial_1^i F(\varepsilon a, z)\|. \end{aligned} \quad (29)$$

This yields the first term on the right of (26). The Lagrange remainder in (28) is estimated independently for the continuum integral over the interval $[a, b]$ and for the discrete integral G_m , in both cases yielding the same contribution to the second term on the right of (26). \square

The third concerns the Gauss formula for the discrete sum $\sum_{j=0}^{N-1} F(x_j)$ on equidistant nodes

$$x_j = -1 + \frac{2j}{N-1}, \quad 0 \leq j \leq N-1 \quad (30)$$

with

$$\frac{2}{N} \sum_{j=0}^{N-1} F(x_j) \approx S_n(F) \equiv \sum_{k=1}^n w_{k,N} F(s_{k,N}), \quad (31)$$

where the quadrature nodes $s_{k,N}$ are the zeros of the so-called Gram polynomial $p_{n,N}$ of degree n . Such polynomials are defined by their orthonormality with respect

to the discrete equidistant sum, namely

$$\sum_{j=0}^{N-1} p_{n,N}(x_j) p_{k,N}(x_j) = \delta_{nk}. \quad (32)$$

The weights $w_{k,N}$ are given by

$$w_{k,N} = \frac{a_{n,N}}{a_{n-1,N}} \frac{2}{N p'_{n,N}(s_{k,N}) p_{n-1,N}(s_{k,N})}, \quad (33)$$

where $a_{n,N}$ denotes the leading coefficient of $p_{n,N}$. For a detailed derivation and discussion, see [1, 2, 19].

Theorem 6. Fix $n \in \mathbb{N}$ such that $n < N$ and let $F: [a, b] \rightarrow X$ be a $2n$ -times differentiable function with $\|F^{(2n)}\| \leq A$ on $[a, b]$. Then

$$\left\| \frac{b-a}{N-1} \sum_{j=0}^{N-1} F(y_j) - \frac{N(b-a)}{2(N-1)} S_n(F) \right\| \leq \frac{16 A (b-a)^{2n+1} n!^4}{(2n+1)(2n)!^3}. \quad (34)$$

Formula $S_n(f)$ is defined with nodes $r_{k,N} = \ell(s_{k,N})$ and the equidistant summation points are given by $y_j = \ell(x_j)$, where ℓ is the affine change of variable (24).

Proof. Assume first that $a = -1$ and $b = 1$; the general case then follows via the affine change of variable ℓ . Assume further that $X = \mathbb{R}$. The general case where X is a complex Banach space follows by applying Lemma 3.

Thus, let H be the unique polynomial of degree $2n-1$ satisfying the Hermite interpolation problem

$$F(s_{k,N}) = H(s_{k,N}), \quad F'(s_{k,N}) = H'(s_{k,N}), \quad k = 1, \dots, n. \quad (35)$$

By Rolle's theorem, for any $x \in [-1, 1]$ there exists $s(x) \in [-1, 1]$ such that

$$F(x) - H(x) = \frac{F^{(2n)}(s)}{(2n)!} q_{n,N}^2(x), \quad (36)$$

where $q_{n,N}$ is the polynomial

$$q_{n,N}(x) = \prod_{k=1}^n (x - x_{k,N}). \quad (37)$$

Since (31) is exact for all polynomials of degree less than $2n-1$,

$$\frac{2}{N} \sum_{j=0}^{N-1} F(x_j) = S_n(H) = S_n(F). \quad (38)$$

Thus, using (36), we estimate

$$\begin{aligned} \left\| \frac{2}{N-1} \sum_{j=0}^{N-1} F(x_j) - \frac{N}{N-1} S_n(F) \right\| &= \frac{2}{N} \sum_{j=0}^{N-1} \|F(x_j) - H(x_j)\| \\ &\leq \frac{2A}{(N-1)(2n)!} \sum_{j=0}^{N-1} q_{n,N}^2(x_j). \end{aligned} \quad (39)$$

Note that $\deg(q_{n,N}) = n$ and $q_{n,N}$ has the same zeros as the Gram polynomial $p_{n,N}$. Hence,

$$p_{n,N} = a_{n,N} q_{n,N}, \quad (40)$$

where the constant $a_{n,N}$ is given by [19]

$$a_{n,N} = \sqrt{\frac{(2N+1)(N-n-1)!}{(N+n)!} \frac{(2n)!(N-1)^n}{2^n n!^2}}. \quad (41)$$

Since $p_{n,N}$ is normalized,

$$\sum_{j=0}^{N-1} q_{n,N}^2(x_j) = \frac{1}{a_{n,N}^2} = \frac{(N+n)!}{(2n+1)(N-n-1)!(N-1)^{2n}} \frac{2^n n!^4}{(2n)!^2}. \quad (42)$$

For $N > n \geq 1$, we have

$$\begin{aligned} \frac{(N+n)!}{(N-n-1)!(N-1)^{2n+1}} &= \frac{(N+n)(N+n-1)\cdots(N-n)}{(N-1)(N-1)\cdots(N-1)} \\ &\leq \left(1 + \frac{1}{n}\right)^{2n+1} \\ &\leq 8, \end{aligned}$$

which completes the proof. \square

In the next section, we will need to approximate a double integral of a function of two variables where one of the integrals is continuous, the other discrete, namely

$$\sum_{j=0}^{N-1} \int_0^1 F(jT, xT, x) dx \quad (43)$$

where $T \approx 1/N \ll 1$. In principle, this is a tensor product construction using Theorem 6 to approximate the discrete sum and Theorem 5 for the continuous integral, except that we need to be careful about uniformity with respect to the small parameter T . To simplify notation later in Section 4, we shall write

$$F(s, Tx, x) \equiv G(s, x; T) \quad (44)$$

and sometimes drop the parametric dependence of G on T for convenience. In the following, we fix $0 < T_0 < \tau_0 < 1$, $0 < \gamma < 1$, and integer $n \geq 1$. Then the assumptions necessary to invoke Theorem 6 resp. Theorem 5 read as follows.

Assumption 1. There exists a constant A such that for every $(x, T) \in [0, 1] \times [0, T_0]$, $s \mapsto G(s, x; T)$ is $2n$ -times differentiable with

$$\sup_{[0, \tau_0] \times [0, 1] \times [0, T_0]} \|\partial_1^{2n} G\| \leq A. \quad (45)$$

Assumption 2. Suppose that

- (i) for every $s \in [0, \tau_0]$ and every $z \in \Sigma_{1/(2\gamma)}(0, 1)$, the map $\zeta \mapsto F(s, \zeta, z)$ is $2n$ -times differentiable on $[0, T_0]$ and there exists a constant A_{an} independent of $s \in [0, \tau_0]$ such that

$$\max_i \sup_{[0, T_0] \times \Sigma_{1/(2\gamma)}(0, 1)} \|\partial_1^i F\| \leq A_{\text{an}}, \quad (46)$$

- (ii) for every $s \in [0, \tau_0]$ and every $\zeta \in [0, T_0]$, the map $z \mapsto \partial_2^j F(s, \zeta, z)$ is analytic on $\Sigma_{1/(2\gamma)}(0, 1)$ for $i = 0, \dots, 2n-1$ and there exists a constant A_{dif} independent of $s \in [0, \tau_0]$ such that

$$\sup_{[0, T_0] \times [0, 1]} \|\partial_2^k F(s, \cdot, \cdot)\| \leq A_{\text{dif}}. \quad (47)$$

The following proposition then combines estimates (34) and (26) in the form required later.

Proposition 7. *Under Assumption 1 and Assumption 2, there exists a constant $C = C(F, \gamma, \tau_0, T_0, n)$ such that for any $(\tau, T) \in (0, \tau_0] \times (0, T_0]$ and $(m, N) \in \mathbb{N}^2$ with $T = \tau/N$ and $n < N$, we have*

$$\left\| T \sum_{j=0}^{N-1} \int_0^1 G(jT, x; T) dx - \frac{\tau}{4} \sum_{i=1}^n \sum_{k=1}^m w_{i,N} \omega_k G(r_{i,N}, \eta_k; T) \right\| \leq C \tau (\gamma^{2m} + \tau^{2n}). \quad (48)$$

Proof. Using Theorem 6, we find that

$$T \sum_{j=0}^{N-1} G(jT, x; T) = \frac{\tau}{2} \sum_{i=1}^n w_{i,N} G(r_{i,N}, x; T) + R(x, T), \quad (49)$$

where, in view of Assumption 1, estimate (34) implies that there exists a constant C depending on $\sup_{[0, \tau_0] \times [0, 1] \times [0, T_0]} \|\partial_1^{2n} G\|$ and n such that

$$\|R(x, T)\| \leq C \tau^{2n+1}. \quad (50)$$

Since $F(r_{i,N}, \cdot, \cdot)$ satisfies Assumption 2, it satisfies the assumption of Theorem 5 on $[0, T_0] \times \Sigma_{1/(2\gamma)}(0, 1)$ for each $r_{i,N}$. Thus, taking the integral of (49) over $[0, 1]$ and using estimate (26), we obtain (48). \square

Remark 2. The main parameters to make the right hand side small are τ and m . The parameter $\gamma \in (0, 1)$, on the other hand, is fixed. We ensure smallness of the right hand side of (48) by first choosing τ sufficiently small, then m sufficiently large such that the first term on the right of (48) is no larger than the second term.

4. UNIFORMLY ACCURATE SCHEMES

Following [4], we introduce “twisted variables” in which the linear operator in the equation is uniform as $c \rightarrow \infty$. The twisting technique was also used in an earlier paper of Castella *et al.* [5] who, in a related context, developed an averaging technique for highly-oscillatory Hamiltonian problems. In a first change of variables, we set

$$U = \phi - c^{-2} B_c^{-1} J \dot{\phi}, \quad (51a)$$

$$V = \phi + c^{-2} B_c^{-1} J \dot{\phi}. \quad (51b)$$

In terms of the variables U and V , equation (1) reads

$$J \dot{U} = -c^2 B_c U + B_c^{-1} f\left(\frac{1}{2}(U + V), t\right), \quad (52a)$$

$$J \dot{V} = c^2 B_c V - B_c^{-1} f\left(\frac{1}{2}(U + V), t\right). \quad (52b)$$

As a second change of variables, we define

$$u = e^{-c^2 t J} U, \quad v = e^{c^2 t J} V. \quad (53)$$

In terms of u and v , system (52) takes the form

$$\dot{u} = J A_c u - J B_c^{-1} e^{-c^2 t J} f\left(\frac{1}{2}(e^{c^2 t J} u + e^{-c^2 t J} v), t\right), \quad (54a)$$

$$\dot{v} = -J A_c v + J B_c^{-1} e^{c^2 t J} f\left(\frac{1}{2}(e^{c^2 t J} u + e^{-c^2 t J} v), t\right). \quad (54b)$$

We can write this system more compactly in terms of the vector-valued functions $W = (U, V)^\top$ and $w = (u, v)^\top$. Letting A_c and B_c act diagonally on $\mathcal{D}(A_c) \times \mathcal{D}(A_c)$ and defining

$$\mathcal{J} = \begin{pmatrix} J & 0 \\ 0 & -J \end{pmatrix}, \quad (55a)$$

$$\mathcal{F}(W, t) = (-J, J)^\top f\left(\frac{1}{2}(U + V), t\right), \quad (55b)$$

we can write

$$\dot{w} = \mathcal{J}A_c w + B_c^{-1} e^{-c^2 t \mathcal{J}} \mathcal{F}(e^{c^2 t \mathcal{J}} w, t). \quad (56)$$

Let $\tau > 0$ be the time step of the numerical scheme. We write $t_i = i\tau$ for $i = 0, 1, 2, \dots$ and apply the Duhamel formula, so that

$$\begin{aligned} w(t_i + \tau) &= e^{\tau \mathcal{J} A_c} w(t_i) \\ &+ B_c^{-1} \int_0^\tau e^{(\tau-s) \mathcal{J} A_c} e^{-c^2(t_i+s) \mathcal{J}} \mathcal{F}(e^{c^2(t_i+s) \mathcal{J}} w(t_i + s), s) ds. \end{aligned} \quad (57)$$

Since we are free to adapt the time τ of what is to emerge as the numerical scheme, it is convenient to select τ as an integer multiple of the fast period $T = 2\pi/c^2$ so that $\tau = NT$ for some $N \in \mathbb{N}$. As $e^{s \mathcal{J}} = \cos(s)I + \sin(s)\mathcal{J}$,

$$e^{\pm c^2 t_i \mathcal{J}} = e^{2\pi i N \mathcal{J}} = I \quad (58)$$

whenever i is integer. Thus, such factors drop out of all expressions further below, reducing the computational cost of the scheme.

The two following assumptions on the nonlinearity f and on the solution of (57) are required for the rigorous analysis of convergence.

Assumption 3. For given $n \in \mathbb{N}$ and $\mathcal{T}_0 > 0$, we assume that f satisfies the following:

- (i) $t \mapsto f(u, t)$ is $2n$ -times differentiable for any $u \in \mathcal{D}(L^{2n})$,
- (ii) $x \mapsto f(e^{2\pi x J} u, t)$ has an analytic extension to $\Sigma_{\frac{1}{2\gamma}}(0, 1)$ for some $\gamma \in (0, 1)$ for any $t \in [0, \mathcal{T}_0]$.
- (iii) f is Lipschitz with respect to the first argument on bounded sets of X with a constant uniform in $t \in [0, \mathcal{T}_0]$.
- (iv) For any $t \in [0, \mathcal{T}_0]$, $f(\cdot, t): \mathcal{D}(L^{2n}) \rightarrow X$ is $2n$ -times Gâteaux differentiable such that $D^k f(u, t) \in \mathcal{L}(\mathcal{D}(L^{2n-\alpha_k}), \mathcal{D}(L^{2n-|\alpha_k|}))$ for every $k = 1, \dots, 2n$, $u \in \mathcal{D}(L^{2n})$, and multi-index $\alpha_k = (j_1, \dots, j_k)$ for which each component is larger than 1 and $|\alpha_k| \leq 2n$.

Here, $\mathcal{D}(L^{2n-\alpha_k})$ refers to the direct product $\mathcal{D}(L^{2n-j_1}) \times \dots \times \mathcal{D}(L^{2n-j_k})$.

Assumption 4. For given n , in the setting of Assumption 3, there exists $\mathcal{T} \in (0, \mathcal{T}_0]$ and $K > 0$ independent of c such that

$$\sup_{0 \leq t \leq \mathcal{T}} \|w(t)\|_{\mathcal{D}(L^n)} \leq K. \quad (59)$$

Remark 3. Assumption 3 includes a wide range of nonlinearities. It is easy to verify that polynomial nonlinearities as well as the nonlinearity of the semilinear Klein–Gordon equation as introduced in Section 1 satisfy this requirement.

Remark 4. To see how the differentiability requirement (iv) arises, consider the following example, which is a simplified version of the estimates which arise in the

analysis of the numerical scheme below. Take $g(s) = e^{sJA_c} f(h(s))$, $h(s) = e^{sJA_c} u$, $u \in \mathcal{D}(L^n)$ and $n = 1$. Since

$$\begin{aligned} e^{-sJA_c} g''(s) &= -A_c^2 f(h(s)) + 2JA_c Df(h(s)) h'(s) \\ &\quad + D^2 f(h(s)) [h'(s), h'(s)] + Df(h(s)) h''(s), \end{aligned} \quad (60)$$

$\|g''\|_X$ is uniformly bounded in c provided

$$Df(u) \in \mathcal{L}(\mathcal{D}(L^{2n-1}), \mathcal{D}(L^{2n-1})) = \mathcal{L}(\mathcal{D}(L)), \quad (61a)$$

$$Df(u) \in \mathcal{L}(\mathcal{D}(L^{2n-2}), \mathcal{D}(L^{2n-2})) = \mathcal{L}(X), \quad (61b)$$

$$D^2 f(u) \in \mathcal{L}(\mathcal{D}(L^{2n-(1,1)}), \mathcal{D}(L^{2n-2})) = \mathcal{L}(\mathcal{D}(L) \times \mathcal{D}(L), X). \quad (61c)$$

This suffices to satisfy Assumption 1 for $G(s, x, T) = g(s)$.

Remark 5. Assumption 4 holds provided the initial data satisfies the bound

$$\|w(0)\|_{\mathcal{D}(L^n)} \leq K_0 \quad (62)$$

where the constant K_0 does not depend on c . This can be proved directly from the Duhamel formula (57) as the operators B_c^{-1} , e^{sJA_c} , and $e^{sc^2\mathcal{J}}$ are bounded uniformly in c in the strong operator topology of $\mathcal{D}(L^n)$.

To guarantee uniform convergence with respect to c , we make the following important observation which effectively asserts that the time derivative \dot{w} is bounded uniformly in c .

Lemma 8. *The solution w of (57) satisfies*

$$\|w(t_i + s) - w(t_i)\| \leq \frac{s}{2} \|w(t_i)\|_{\mathcal{D}(L)} + s \sup_{\sigma \in [0, s]} \|\mathcal{F}(e^{c^2(t_i + \xi)\mathcal{J}} w(t_i + \sigma))\|. \quad (63)$$

Proof. The proof is a direct application of estimate (vi) in Lemma 2 and the fact that $\|B_c^{-1}\| \leq 1$. \square

In a first step, we define a sequence of “pre-schemes” $\Phi_l: X \times \mathbb{R} \rightarrow X$ which provide consistent approximations to the right hand side of the Duhamel formula (57) to order τ^{l+1} , namely

$$\Phi_1(w, z) = e^{zJA_c} w - B_c^{-1} \int_0^z e^{-c^2 s \mathcal{J}} \mathcal{F}(e^{c^2 s \mathcal{J}} w, s) ds, \quad (64a)$$

$$\Phi_{l+1}(w, z) = e^{zJA_c} w - B_c^{-1} \int_0^z e^{(z-s)JA_c} e^{-c^2 s \mathcal{J}} \mathcal{F}(e^{c^2 s \mathcal{J}} \Phi_l(w, s), s) ds. \quad (64b)$$

The pre-schemes approximate the true solution in the following sense.

Lemma 9. *Under Assumption 3 (iii), let w be a solution for (57) satisfying Assumption 4 for $n = 1$, and fix $l \in \mathbb{N}^*$. Then there exist constants C_l independent of c such that all $s \geq 0$,*

$$\|w(t_i + s) - \Phi_l(w(t_i), s)\| \leq C_l s^{l+1}. \quad (65)$$

Proof. We set $R_l(w(t_i), s) = w(t_i + s) - \Phi_l(w(t_i), s)$ and proceed by induction. When $l = 1$,

$$\begin{aligned} R_1(w(t_i), s) &= B_c^{-1} \int_0^s e^{-c^2 \sigma \mathcal{J}} \mathcal{F}(e^{c^2 \sigma \mathcal{J}} w(t_i), \sigma) d\sigma \\ &\quad - B_c^{-1} \int_0^s e^{(s-\sigma)JA_c} e^{-c^2 \sigma \mathcal{J}} \mathcal{F}(e^{c^2 \sigma \mathcal{J}} w(t_i + \sigma), \sigma) d\sigma. \end{aligned} \quad (66)$$

The estimate on R_1 follows by using Lemma 8 to freeze $w(t_i + \sigma)$ and Lemma 2(iv) to remove the operator $e^{(s-\sigma)\mathcal{J}A_c}$ in the second integral in (66). For $l \geq 1$,

$$\begin{aligned} R_{l+1}(w(t_i), s) &= B_c^{-1} \int_0^s e^{(s-\sigma)\mathcal{J}A_c} e^{-c^2\sigma\mathcal{J}} \mathcal{F}(e^{c^2\sigma\mathcal{J}} \Phi_l(w(t_i), \sigma), \sigma) d\sigma \\ &\quad - B_c^{-1} \int_0^s e^{(s-\sigma)\mathcal{J}A_c} e^{-c^2\sigma\mathcal{J}} \mathcal{F}(e^{c^2\sigma\mathcal{J}} w(t_i + \sigma), \sigma) d\sigma. \end{aligned} \quad (67)$$

By Lemma 2 and the fact that f is Lipschitz on bounded sets of X with respect to the first argument, there exists a constant C independent of c such that

$$\|R_{l+1}(w(t_i), s)\| \leq C s \sup_{\sigma \leq s} \|R_l(w(t_i), \sigma)\|. \quad (68)$$

This completes the proof. \square

While the operator A_c and the associated semi-group $e^{t\mathcal{J}A_c}$ are uniformly well-behaved as $c \rightarrow \infty$, the integrals in (64) still contain highly oscillatory terms with a *single* fast frequency. For the latter, effective numerical quadrature is possible [19]. Following the strategy developed there, we split $z/T \equiv N_z + \theta_z$ into its integer part $N_z = \lfloor z/T \rfloor$ and fractional part $\theta_z = z/T - N_z$. Then the integral in (64a) can be written

$$\begin{aligned} &B_c^{-1} \int_0^z e^{-c^2s\mathcal{J}} \mathcal{F}(e^{c^2s\mathcal{J}} w, s) ds \\ &= B_c^{-1} \sum_{j=0}^{N_z-1} \int_{jT}^{(j+1)T} e^{-c^2s\mathcal{J}} \mathcal{F}(e^{c^2s\mathcal{J}} w, s) ds \\ &\quad + B_c^{-1} \int_{N_z T}^z e^{-c^2s\mathcal{J}} \mathcal{F}(e^{c^2s\mathcal{J}} w, s) ds \\ &= T \sum_{j=0}^{N_z-1} \int_0^1 G_0(jT, \sigma) d\sigma + T \int_0^{\theta_z} G_0(N_z T, \sigma) d\sigma \end{aligned} \quad (69)$$

with

$$G_0(\rho, \sigma) = B_c^{-1} e^{-2\pi\sigma\mathcal{J}} \mathcal{F}(e^{2\pi\sigma\mathcal{J}} w, \rho + \sigma T) \quad (70)$$

and where, in the second equality of (69), we have used (58). Analogously, the integral in (64b) can be written

$$\begin{aligned} &B_c^{-1} \int_0^z e^{-s\mathcal{J}A_c} e^{-c^2s\mathcal{J}} \mathcal{F}(e^{c^2s\mathcal{J}} \Phi_l(w, s), s) ds \\ &= T \sum_{j=0}^{N_z-1} \int_0^1 G[\Phi_l](jT, \sigma) d\sigma + T \int_0^{\theta_z} G[\Phi_l](N_z T, \sigma) d\sigma, \end{aligned} \quad (71)$$

where, for $\Upsilon: X \times \mathbb{R} \rightarrow X$,

$$G[\Upsilon](\rho, \sigma) = B_c^{-1} e^{-(\rho+\sigma T)\mathcal{J}A_c} e^{-2\pi\sigma\mathcal{J}} \mathcal{F}(e^{2\pi\sigma\mathcal{J}} \Upsilon(w, \rho + \sigma T), \rho + \sigma T). \quad (72)$$

Altogether, (64) then takes the form

$$\Phi_1(w, z) = e^{z\mathcal{J}A_c} w - T \sum_{j=0}^{N_z-1} \int_0^1 G_0(jT, \sigma) d\sigma - T \int_0^{\theta_z} G_0(N_z T, \sigma) d\sigma, \quad (73a)$$

$$\Phi_{l+1}(w, z) = e^{z\mathcal{J}A_c} \left(w - T \sum_{j=0}^{N_z-1} \int_0^1 G[\Phi_l](jT, \sigma) d\sigma - T \int_0^{\theta_z} G[\Phi_l](N_z T, \sigma) d\sigma \right). \quad (73b)$$

We now use the approximate the integrals in (73) by classical Gauss quadrature and the sums by Gauss summation to obtain

$$\begin{aligned} \Psi_1(w, z) &= e^{z\mathcal{J}A_c} w - \frac{N_z T}{4} \sum_{j=1}^n \sum_{k=0}^m w_{j, N_z} \omega_k G_0(r_{j, N_z}, \eta_k) \\ &\quad - \frac{\theta_z T}{2} \sum_{k=0}^m \omega_k G_0(N_z T, \theta_z \eta_k), \end{aligned} \quad (74a)$$

$$\begin{aligned} \Psi_{l+1}(w, z) &= e^{z\mathcal{J}A_c} \left(w - \frac{N_z T}{4} \sum_{j=1}^n \sum_{k=0}^m w_{j, N_z} \omega_k G[\Psi_l](r_{j, N_z}, \eta_k) \right. \\ &\quad \left. - \frac{\theta_z T}{2} \sum_{k=0}^m \omega_k G[\Psi_l](N_z T, \theta_z \eta_k) \right). \end{aligned} \quad (74b)$$

Remark 6. For a scheme of global order l , we use Ψ_l with $z = \tau$ as the time stepper. At the top level, the second sum in (74a) or (74b) does not contribute. However, when $l \geq 2$, the inner evaluations of $\Psi_{l-1}, \Psi_{l-2}, \dots$ will generally be evaluated at points z that are not integer multiples of T , so that their z -arguments have to be re-split into the respective integer (N_z) and fractional (θ_z) multiples of T . Thus, in general, the second sum on the right of (74) is required for consistency.

Remark 7. Note that in the case where \mathcal{F} is constant with respect to the second variable, the function $G_0 = G_0(x)$ is one-variable periodic function. Thus, the approximation from Theorem 4 can also be used to define a first order scheme ($l = 1$) with accuracy that is exponential in the number of nodes. More specifically, for $\tau = NT$,

$$\begin{aligned} \Phi_1(w, \tau) &= e^{\tau\mathcal{J}A_c} w - \tau \int_0^1 G_0(x) dx \\ &= e^{\tau\mathcal{J}A_c} w - \frac{\tau}{m} \sum_{k=0}^{m-1} G_0\left(\frac{k}{m}\right) + \mathcal{O}(\tau e^{-dm}) \end{aligned} \quad (75)$$

for some $d > 0$.

Lemma 10. *Let $l, n \in \mathbb{N}^*$ and $w \in \mathcal{D}(L^{2n})$. Fix $0 < z_0 < 1$, $c_0 > 0$, and assume that f satisfies Assumption 3, with analyticity property (ii) valid on the ellipse $\Sigma_{1/(2\gamma)}(0, 1)$ for some $\gamma \in (0, 1)$. Then there exists $C_l = C_l(f, \|w\|_{\mathcal{D}(L^{2n})}, c_0, z_0, n)$ such that for all $m \in \mathbb{N}^*$ and $z \leq z_0 < 1$,*

$$\|\Psi_l(w, z) - \Phi_l(w, z)\| \leq C_l z (z^{2n} + \gamma^{2m}). \quad (76)$$

Proof. In view of the expression for each G_0 and G , there exist F_0 and F such that

$$G_0(\rho, \gamma, T) = F_0(\rho, T\sigma, \sigma), \quad G(\rho, \gamma, T) = F(\rho, T\sigma, \sigma), \quad (77)$$

where, since $w \in \mathcal{D}(L^{2n})$ and f satisfies Assumption 3, F_0 and F satisfy the conditions of Proposition 7 on $[0, z_0] \times [0, 2\pi/c_0^2] \times \Sigma_{1/(2\gamma)}(0, 1)$.

We set $S_l(w, z) = \Psi_l(w, z) - \Phi_l(w, z)$ and proceed by induction. For $l = 1$, we can directly use Proposition 7 for the difference of first terms and Theorem 5 for the difference of second terms, (76) holds true as stated. For $l > 1$, we have

$$\begin{aligned}
& e^{-z\mathcal{J}A_c} S_{l+1}(w, z) \\
&= -\frac{N_z T}{4} \sum_{j=1}^n \sum_{k=0}^m w_{j, N_z} \omega_k G[\Psi_l](r_{j, N_z}, \eta_k) - \frac{\theta_z T}{2} \sum_{k=0}^m \omega_k G[\Psi_l](N_z T, \theta_z \eta_k) \\
&+ \frac{N_z T}{4} \sum_{j=1}^n \sum_{k=0}^m w_{j, N_z} \omega_k G[\Phi_l](r_{j, N_z}, \eta_k) + \frac{\theta_z T}{2} \sum_{k=0}^m \omega_k G[\Phi_l](N_z T, \theta_z \eta_k) \\
&- \frac{N_z T}{4} \sum_{j=1}^n \sum_{k=0}^m w_{j, N_z} \omega_k G[\Phi_l](r_{j, N_z}, \eta_k) - \frac{\theta_z T}{2} \sum_{k=0}^m \omega_k G[\Phi_l](N_z T, \theta_z \eta_k) \\
&+ T \sum_{j=0}^{N_z-1} \int_0^1 G[\Phi_l](jT, \sigma) d\sigma + T \int_0^{\theta_z} G[\Phi_l](N_z T, \sigma) d\sigma
\end{aligned} \tag{78}$$

We write $S_{l+1}^{(1)}(w, z)$ and $S_{l+1}^{(2)}(w, z)$ to denote the first two and the last two lines on the right of (78), respectively. As f is Lipschitz with respect to the first argument on bounded sets of X , there exist K_1 and K_2 , each depending on f , such that

$$\begin{aligned}
\|S_{l+1}^{(1)}(w, z)\| &\leq K_1 z \sup_{j,k} \|S_l(w, r_{j, N_z} + \eta_k T)\| + K_2 z \sup_k \|S_l(w, N_z T + \theta_z \eta_k T)\| \\
&\leq (K_1 + K_2) C_l z^2 (z^{2n} + \gamma^{2m}).
\end{aligned} \tag{79}$$

On the other hand, using Proposition 7, there exists a constant K_3 depending on f , $\|w\|_{\mathcal{D}(L^{2n})}$, c_0 , z_0 , and n such that

$$\|S_{l+1}^{(2)}(w, z)\| \leq K_3 z (z^{2n} + \gamma^{2m}). \tag{80}$$

Thus, combining (79) and (80), we conclude that there exists a constant C_{l+1} depending on f , $\|w\|_{\mathcal{D}(L^{2n})}$, c_0 , z_0 and n such that

$$\|S_{l+1}(w, z)\| \leq C_{l+1} z (z^{2n} + \gamma^{2m}), \tag{81}$$

which concludes the proof. \square

As stated before, we select the time step τ to be an integer multiple of the fast period T so that $\tau = NT$ for some $N \in \mathbb{N}$. As a numerical approximation to the exact solution w at time t_{i+1} , we take the scheme

$$w_{i+1} = \Psi_l(w_i, \tau), \tag{82a}$$

$$w_0 = \begin{pmatrix} \phi_0 \\ \phi_0 \end{pmatrix} - c^{-2} \mathcal{J} B_c^{-1} \begin{pmatrix} \phi_0' \\ \phi_0' \end{pmatrix}. \tag{82b}$$

5. CONVERGENCE ANALYSIS

The scheme (82) satisfies the following global estimate.

Theorem 11. *Let f satisfies Assumption 3, with analyticity property (i) valid on the ellipse $\Sigma_{1/(2\gamma)}(0, 1)$ for some $\gamma \in (0, 1)$. Fix $l \in \mathbb{N}^*$, $c_0 > 0$, and let $n = \lfloor \frac{l+1}{2} \rfloor$. Assume further that there exists $\mathcal{K} > 0$ such that for every $c \geq c_0$,*

$$\|\phi_0\|_{\mathcal{D}(L^{2n})} + c^{-2} \|B_c^{-1}\phi'_0\|_{\mathcal{D}(L^{2n})} \leq \mathcal{K}. \quad (83)$$

Then there exist $\mathcal{T} > 0$ and $C = C(f, \mathcal{K}, \mathcal{T}, c_0, n)$ such that for all $c \geq c_0$, $\tau \in \frac{2\pi}{c^2}\mathbb{N}$, $t_i \leq \mathcal{T}$, and $m \in \mathbb{N}^$,*

$$\|\phi_i - \phi(t_i)\| \leq C(\tau^l + \gamma^{2m}) \quad (84)$$

where ϕ solves (1) and

$$\phi_i = \frac{(w_i)_1 + (w_i)_2}{2} \quad (85)$$

with w_i given by (82).

Proof. Note first that for every $c \geq c_0$,

$$\|w_0\|_{\mathcal{D}(L^{2n})} \leq \|\phi_0\|_{\mathcal{D}(L^{2n})} + \|c^{-2}JB_c^{-1}\phi'_0\|_{\mathcal{D}(L^{2n})} \leq \mathcal{K}. \quad (86)$$

Thus, there exist two constants $\mathcal{T}, K > 0$ depending on c_0 and w_0 for which Assumption 4 is satisfied. Lemmas 10 and 9 allow us to write

$$\begin{aligned} w(t_i + \tau) &= \Phi_l(w(t_i), \tau) + R_l(w(t_i), \tau) \\ &= \Psi_l(w(t_i), \tau) + R_l(w(t_i), \tau) - S_l(w(t_i), \tau). \end{aligned} \quad (87)$$

Setting $e_i = \|w(t_i) - w_i\|$, we now split the error as follows:

$$\begin{aligned} e_{i+1} &\leq \|R_l(w(t_i), \tau)\| + \|S_l(w(t_i), \tau)\| \\ &\quad + \|\Psi_l(w(t_i), \tau) - \Psi_l(w_i, \tau)\|. \end{aligned} \quad (88)$$

Recalling that \mathcal{F} is Lipschitz on X and arguing by induction on l , we find that there exists a constant $C_1 > 0$ depending on f such that

$$\|\Psi_l(w(t_i), \tau) - \Psi_l(w_i, \tau)\| \leq (1 + C_1\tau)^l e_i. \quad (89)$$

By Lemma 9 and 10, there exists a constant $C_2 > 0$ depending on $f, \mathcal{K}, \mathcal{T}, c_0$, and n such that

$$\|R_l(w(t_i), \tau)\| + \|S_l(w(t_i), \tau)\| \leq C_2\tau(\tau^l + \gamma^{2m}). \quad (90)$$

Then, (88) reads

$$e_{i+1} \leq (1 + C_1\tau)^l e_i + C_2\tau(\tau^l + \gamma^{2m}). \quad (91)$$

Thus, we find by induction that

$$e_i \leq (1 + C_1\tau)^{il} e_0 + C_2 \frac{(1 + C_1\tau)^{il} - 1}{C_1} (\tau^l + \gamma^{2m}). \quad (92)$$

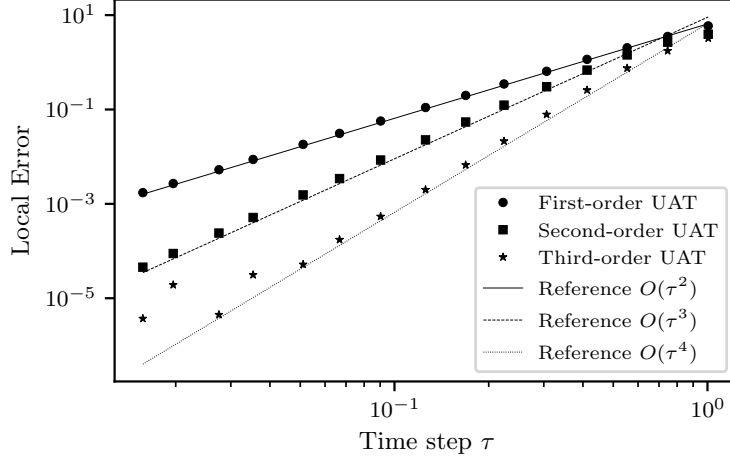
Since $e_0 = 0$ and $1 + x \leq e^x$, we obtain

$$e_i \leq C_2 \frac{e^{C_1 l \tau} - 1}{C_1} (\tau^l + \gamma^{2m}) \equiv C(\tau^l + \gamma^{2m}). \quad (93)$$

To obtain the final estimate, we undo the variable twist, noting that

$$\phi(t_i) = \frac{(e^{c^2 t_i J}(w(t_i)))_1 + (e^{-c^2 t_i J}(w(t_i)))_2}{2} = \frac{w(t_i)_1 + w(t_i)_2}{2}. \quad (94)$$

Then (93) directly implies estimate (84). \square

FIGURE 1. Scaling of the local error with the time step τ .

6. NUMERICAL TESTS

We now demonstrate the scaling behavior of the new uniformly accurate time integrators (UAT) in a simple test case where an explicit reference solution is available. Our example has $\phi: [0, T] \times \mathbb{T}^d \rightarrow \mathbb{C}$ with $L = \delta - \Delta$ and $f(\phi) = |\phi|^2\phi$ with some $\delta > 0$. Then, for arbitrary $a \in \mathbb{R}^d$, the function

$$\phi(t, x) = \sqrt{\delta + |a|^2} e^{i(ct+ax)} \quad (95)$$

is a solution of (1) with

$$\phi_0 = \sqrt{\delta + |a|^2} e^{ia \cdot x}, \quad \phi'_0 = ic\sqrt{\delta + |a|^2} e^{ia \cdot x}. \quad (96)$$

For simplicity, we consider only solutions with no dependence on x , i.e. $a = 0$, where

$$\phi(t, x) = \phi(t) = \sqrt{\delta} e^{ict}. \quad (97)$$

The order of the Gaussian quadrature approximating the inner integral in (73) is chosen as $m = 6, 8, 10$ at level $l = 0, 1, 2$. Theoretically, in view of (84), m should be chosen so that

$$m \approx \frac{\ln(\tau)}{2 \ln(\gamma)} l. \quad (98)$$

However, as we do not have any access to a good estimate for γ , we determined a minimal choice of m empirically.

In Figure 1, we confirm numerically the theoretical convergence rate with respect to τ for the first, second and third order schemes given by (82). Shown is the local error

$$E_{\text{loc}}(\tau) = \|\phi_1 - \phi(\tau)\|, \quad (99)$$

which corresponds to $i = 1$ in (84), for fixed $c = 200$ as the time step τ is varied. As explained in Section 4, we work with time steps that are integer multiples of the fast period, i.e., $\tau = \frac{2\pi}{c^2} k$ for integer k . It is possible to modify the code such that arbitrary time steps are possible. However, this would require retaining all factors $e^{\pm c^2 t_i}$ in the generating formula (57) and all expressions that follow, and the second

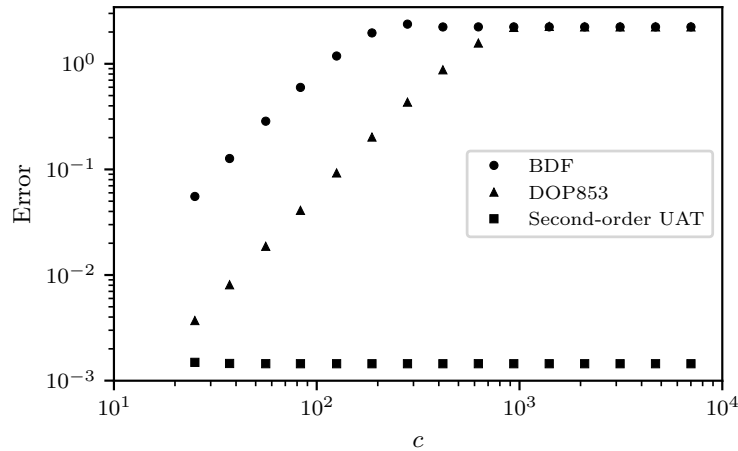


FIGURE 2. Comparison of the second-order uniformly accurate time integrator with two of the standard solvers from Scipy: the implicit variable-order backward differentiation scheme and the explicit Dormand–Prince embedded order 8(5,3) Runge–Kutta scheme.

sum in (74) would already appear at the top level of the recursion, cf. Remark 6. As there is no advantage of doing so, we did not implement this general case.

Figure 1 shows, in particular, that the local error of the third order method scales like τ^4 , thus the global error will scale like τ^3 , except for rather small values of τ where the limitations of double-precision floating point begin to matter. In general, floating-point errors might occur when we use the quadrature formula on very small intervals. As Figure 1 shows, this occurs when using the third order scheme for small values of τ , since calling the function $\Psi_3(w, \tau)$ defined in (74b) includes implicit calls of $\Psi_2(w, y)$ and $\Psi_1(w, x)$ for x, y with $x \ll y \ll \tau$, which means that the lengths of the subintervals on which we use the Gauss formula for the discrete sums when calling $\Psi_2(w, y)$ and $\Psi_1(w, x)$ are very small.

Figure 2 illustrates the uniformity of the error as a function of c using the second order UAT integrator compared with the explicit Dormand–Prince embedded order 8(5,3) Runge–Kutta (“DOP853”) scheme [17] and the implicit multi-step variable-order (1 to 5) method based on a backward differentiation formula (“BDF”) for the derivative approximation [13]. Error performance of the new scheme is uniform, while the error increases with c for the built-in schemes, indicating that their error indicator heuristics are insufficient for dealing with such extreme multi-scale dynamics. We remark that the coefficients of the Gauss summation formula need to be recomputed for every value of N . The main step of this computation is finding the roots of a polynomial of degree n , with N -dependent coefficients. For practical, small values of n , the associated cost is not a significant part of the overall cost of computation.

In Figure 3, we compare the computation time for single time step using the second order integrator with the solvers used in Figure 2. We see that computation time goes up quadratically in c , as expected, while computation time of the UAT

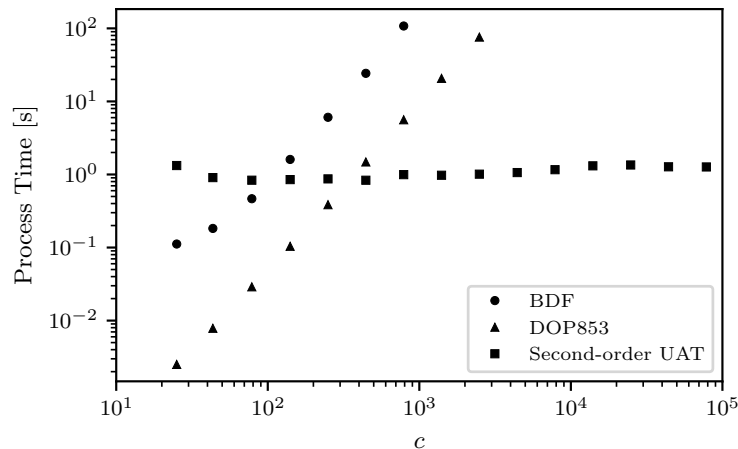


FIGURE 3. Computation time for a single time step of the uniformly accurate time integrator vs. the computation time over the same interval of time using the built-in solvers, controlling their error tolerances so that their error is no more than 30% larger than the error of the uniform scheme.

scheme is constant by design. Timings refer to our reference implementation in Python, which is available as supplementary material, on a single core of an Intel i7 mobile processor, without any attempt at speed-optimizing the code which is bottlenecked in the Python interpreter for this low-dimensional test problem. A more involved study on approximate slow manifolds for semilinear equations of Klein–Gordon type is current work-in-progress and will be reported on separately.

ACKNOWLEDGMENTS

The work was supported by German Research Foundation (DFG) grants OL-155/6-2 and MO-4162/1-1. The authors acknowledge additional support through German Research Foundation Collaborative Research Center TRR 181 under project number 274762653.

REFERENCES

- [1] I. AREA, D. K. DIMITROV, E. GODOY, AND V. G. PASCHOA, *Approximate calculation of sums I: Bounds for the zeros of Gram polynomials*, SIAM J. Numer. Anal., 52 (2014), pp. 1867–1886.
- [2] ———, *Approximate calculation of sums II: Gaussian type quadrature*, SIAM J. Numer. Anal., 54 (2016), pp. 2210–2227.
- [3] W. BAO, Y. CAI, AND X. ZHAO, *A uniformly accurate multiscale time integrator pseudospectral method for the Klein-Gordon equation in the nonrelativistic limit regime*, SIAM J. Numer. Anal., 52 (2014), pp. 2488–2511.
- [4] 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 (2018), pp. 1227–1254.
- [5] F. CASTELLA, P. CHARTIER, AND E. FAOU, *An averaging technique for highly oscillatory Hamiltonian problems*, SIAM J. Numer. Anal., 47 (2009), pp. 2808–2837.

- [6] P. CHARTIER, M. LEMOU, AND F. MÉHATS, *Highly-oscillatory evolution equations with multiple frequencies: averaging and numerics*, Numer. Math., 136 (2017), pp. 907–939.
- [7] P. CHARTIER, M. LEMOU, F. MÉHATS, AND G. VILMART, *A new class of uniformly accurate numerical schemes for highly oscillatory evolution equations*, Found. Comput. Math., 20 (2020), pp. 1–33.
- [8] M. M. CHAWLA AND M. K. JAIN, *Error estimates for Gauss quadrature formulas for analytic functions*, Math. Comp., 22 (1968), pp. 82–90.
- [9] D. COHEN, E. HAIRER, AND C. LUBICH, *Modulated Fourier expansions of highly oscillatory differential equations*, Found. Comput. Math., 3 (2003), pp. 327–345.
- [10] C. J. COTTER AND S. REICH, *Semigeostrophic particle motion and exponentially accurate normal forms*, Multiscale Model. Simul., 5 (2006), pp. 476–496.
- [11] E. FAOU AND K. SCHRATZ, *Asymptotic preserving schemes for the Klein–Gordon equation in the non-relativistic limit regime*, Numer. Math., 126 (2014), pp. 441–469.
- [12] L. GAUCKLER, *Error analysis of trigonometric integrators for semilinear wave equations*, SIAM J. Numer. Anal., 53 (2015), pp. 1082–1106.
- [13] C. W. GEAR, *The numerical integration of ordinary differential equations*, Mathematics of Computation, 21 (1967), pp. 146–156.
- [14] G. A. GOTTWALD, H. MOHAMAD, AND M. OLIVER, *Optimal balance via adiabatic invariance of approximate slow manifolds*, Multiscale Model. Simul., 15 (2017), pp. 1404–1422.
- [15] G. A. GOTTWALD AND M. OLIVER, *Slow dynamics via degenerate variational asymptotics*, Proc. R. Soc. Lond. Ser. A Math. Phys. Eng. Sci., 470 (2014), p. 20140460.
- [16] E. HAIRER, C. LUBICH, AND G. WANNER, *Geometric numerical integration*, vol. 31 of Springer Series in Computational Mathematics, Springer-Verlag, Berlin, second ed., 2006. Structure-preserving algorithms for ordinary differential equations.
- [17] E. HAIRER, S. P. NORSETT, AND G. WANNER, *Solving Ordinary Differential Equations I: Nonstiff Problems*, Springer Berlin, Heidelberg, second ed., 1993.
- [18] N. MASMOUDI AND K. NAKANISHI, *From nonlinear Klein-Gordon equation to a system of coupled nonlinear Schrödinger equations*, Math. Ann., 324 (2002), pp. 359–389.
- [19] H. MOHAMAD AND M. OLIVER, *Numerical integration of functions of a rapidly rotating phase*, SIAM J. Num. Anal., 59 (2021), pp. 2310–2319.
- [20] K. SCHMÜDGEN, *Unbounded self-adjoint operators on Hilbert Space*, Springer Netherlands, first ed., 2012.
- [21] W. STRAUSS AND L. VAZQUEZ, *Numerical solution of a nonlinear Klein-Gordon equation*, J. Comput. Phys., 28 (1978), pp. 271–278.
- [22] L. N. TREFETHEN AND J. A. C. WEIDEMAN, *The exponentially convergent trapezoidal rule*, SIAM Rev., 56 (2014), pp. 385–458.

(H. Mohamad and M. Oliver) MATHEMATICAL INSTITUTE FOR MACHINE LEARNING AND DATA SCIENCE, KU EICHSTÄTT–INGOLSTADT, 85049 INGOLSTADT, GERMANY

(M. Oliver) CONSTRUCTOR UNIVERSITY, 28759 BREMEN, GERMANY