# Moving forward with the help of backward error analysis: improving symplectic and other numerical methods

Farkas Miklós Seminar on Applied Analysis, BME

### Donát M. Takács

Budapest University of Technology and Economics, Faculty of Mechanical Engineering, Department of Energy Engineering and

HUN-REN Centre for Energy Research, Fusion Plasma Physics Department

takacs@energia.bme.hu

2025. 11. 13.

### Overview

- 1. An introduction to backward error analysis (BEA) for ODEs
- 2. BEA-based compensation I. Modifying system parameters
- 3. BEA-based compensation II. Choosing a better coordinate system

Backward error analysis for ODEs

### Notation, terminology

Autonomous, *d*-dimensional ODE (system), initial value problem:

$$\dot{\mathbf{y}}(t) = \mathbf{f}(\mathbf{y}(t)), \qquad \mathbf{y}(0) = \mathbf{y_0},$$
 (1)

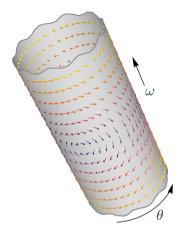
where  $\mathbf{y}: \mathbb{R} \to \mathbb{R}^d$  is the (unknown) solution,  $\mathbf{f}: \mathbb{R}^d \to \mathbb{R}^d$  is the (sufficiently smooth) generating vector field (the "system" being simulated), and  $\mathbf{y_0} \in \mathbb{R}^d$  is the initial condition.

Explicit, one-step numerical method on a uniform grid with time step  $\Delta t$ :

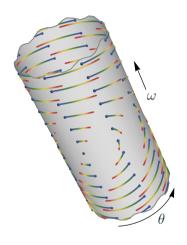
$$\mathbf{y}^{j+1} = \Phi_{\Delta t}(\mathbf{y}^j), \quad j = 0, 1, 2, \dots$$
 (2)

where  $t^j:=j\Delta t$ , with  $j=0,1,2,\ldots$  as the index of the time step,  $\mathbf{y}^j$  is the approximate solution at  $t^j$ . Designed to solve (1) approximately.

Flow of the vector field (integral curves at all points in state space, parametrised by t – a one-parameter map):  $\varphi_t : \mathbb{R} \times \mathbb{R}^d \to \mathbb{R}^d$ 



Vector field  ${\bf f}$ 

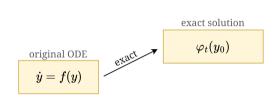


Flow  $\varphi_t$ 

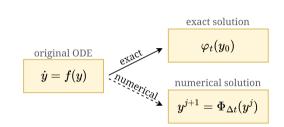
original ODE

$$\dot{y}=f(y)$$

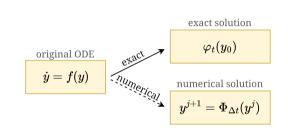
 The flow corresponds to the exact solution of the original problem



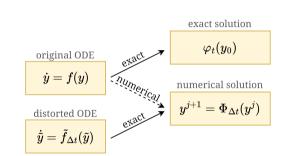
- The flow corresponds to the exact solution of the original problem
- The numerical solution is only approximate



- The flow corresponds to the exact solution of the original problem
- The numerical solution is only approximate
- What is the numerical solution an exact solution of?

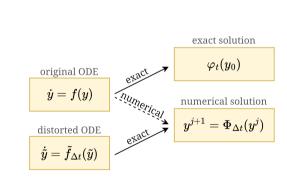


- The flow corresponds to the exact solution of the original problem
- The numerical solution is only approximate
- What is the numerical solution an exact solution of?
- BEA: constructing the distorted ODE



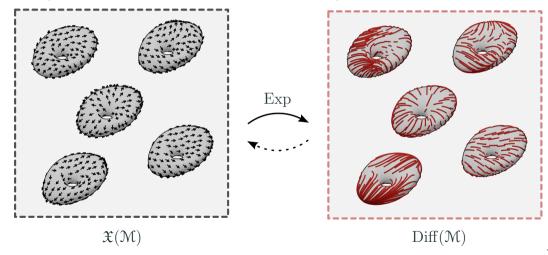
### Backward error analysis for ODEs

- Backward error analysis: treating the approximate solution of  ${\bf f}$  by  $\Phi_{\Delta t}$  as an exact solution of a certain system nearby to the original one, i.e.  $\tilde{{\bf y}}(t^j)={\bf y}^j.$
- This other system is described by its so-called modified equations or distorted equations.
- Corresponding distorted vector field (DVF):  $\tilde{\mathbf{f}}$
- f is an underlying, continuous-time counterpart of the discrete-time numerical method



### Correspondence between the Lie algebra and the Lie group

Lie algebra of vector fields with Lie derivative  $\stackrel{?}{\Leftrightarrow}$  Lie group of smooth maps?



7/38

The distorted vector field (DVF)  $\tilde{\mathbf{f}}$ :

• ...Does it exist?

The distorted vector field (DVF)  $\tilde{\mathbf{f}}$ :

- ...Does it exist?
- usually, it does [1], but...

The distorted vector field (DVF)  $\tilde{\mathbf{f}}$ :

- ...Does it exist?
- usually, it does [1], but...
- ullet it is not autonomous:  $ilde{\mathbf{f}}(\mathbf{y},t)$

The distorted vector field (DVF)  $\tilde{\mathbf{f}}$ :

- ...Does it exist?
- usually, it does [1], but...
- it is not autonomous:  $\tilde{\mathbf{f}}(\mathbf{y},t)$
- ullet or, it can be approximated via an asymptotic series of an autonomous  $ilde{\mathbf{f}}(\mathbf{y})$  as:

$$\tilde{\mathbf{f}}(\tilde{\mathbf{y}}) = \mathbf{f}(\tilde{\mathbf{y}}) + \Delta t \, \mathbf{f}_1(\tilde{\mathbf{y}}) + \Delta t^2 \, \mathbf{f}_2(\tilde{\mathbf{y}}) + \Delta t^3 \, \mathbf{f}_3(\tilde{\mathbf{y}}) + \cdots$$
(3)

The distorted vector field (DVF)  $\tilde{\mathbf{f}}$ :

- ...Does it exist?
- usually, it does [1], but...
- it is not autonomous:  $\tilde{\mathbf{f}}(\mathbf{y},t)$
- ullet or, it can be approximated via an asymptotic series of an autonomous  $ilde{\mathbf{f}}(\mathbf{y})$  as:

$$\tilde{\mathbf{f}}(\tilde{\mathbf{y}}) = \mathbf{f}(\tilde{\mathbf{y}}) + \Delta t \, \mathbf{f}_1(\tilde{\mathbf{y}}) + \Delta t^2 \, \mathbf{f}_2(\tilde{\mathbf{y}}) + \Delta t^3 \, \mathbf{f}_3(\tilde{\mathbf{y}}) + \cdots$$
(3)

Truncated version is "good enough" for "long enough"

### Nice properties of the DVF

$$\tilde{\mathbf{f}}(\tilde{\mathbf{y}}) = \mathbf{f}(\tilde{\mathbf{y}}) + \Delta t \, \mathbf{f}_1(\tilde{\mathbf{y}}) + \Delta t^2 \, \mathbf{f}_2(\tilde{\mathbf{y}}) + \Delta t^3 \, \mathbf{f}_3(\tilde{\mathbf{y}}) + \cdots$$
(3)

Nice properties of the DVF include:

- If  $\Phi_{\Delta t}$  has order p, then  $\mathbf{f}_j(\mathbf{y}) \equiv \mathbf{0}$  for  $j = 1, \dots, p-1$
- ullet If  $\Phi_{\Delta t}$  is a symmetric method, then  $\mathbf{f}_j(\mathbf{y}) \equiv \mathbf{0}$  for all odd j
- Transfer of structure-preserving properties:
  - If  $\Phi_{\Delta t}$  exactly conserves a first integral  $I(\mathbf{y})$ , then the distorted equation also has  $I(\mathbf{y})$  as a first integral
  - If  $\Phi_{\Delta t}$  is symplectic when applied to a Hamiltonian system of the form  $\mathbf{f} = \mathcal{J}_{\mathbf{c}} \, \mathbf{D} H(\mathbf{y})$ , then the distorted equation is also Hamiltonian  $\rightarrow$  with a distorted Hamiltonian:  $\tilde{\mathbf{f}} = \mathcal{J}_{\mathbf{c}} \, \mathbf{D} \tilde{H}(\mathbf{y})$  Central result for symplectic methods: [2]
  - etc.

Fruitful for the analysis of structure-preserving numerical methods

# Example

Example: explicit Euler (EE) method, system:  $\dot{\mathbf{y}} = \mathbf{f}(\mathbf{y}(t))$ 

$$\mathbf{y}^{j+1} = \mathbf{y}^j + \Delta t \, \mathbf{f}(\mathbf{y}^j) =: \Phi_{\Delta t}(\mathbf{y}^j)$$

Looking for the DVF in the form (ansatz):

$$\tilde{\mathbf{f}}(\tilde{\mathbf{y}}) = \mathbf{f}(\tilde{\mathbf{y}}) + \Delta t \, \mathbf{f}_1(\tilde{\mathbf{y}}) + \Delta t^2 \, \mathbf{f}_2(\tilde{\mathbf{y}}) + \cdots$$

Taylor expansion of the distorted solution  $\tilde{\mathbf{v}}$  around t:

$$\mathbf{y}(t + \Delta t) = \mathbf{y}(t)$$

a 
$$\tilde{\mathbf{f}}(\tilde{\mathbf{v}})$$
 (chain rule):

Expressing this via 
$$\tilde{\mathbf{f}}(\tilde{\mathbf{y}})$$
 (chain rule):

$$_{\Delta t}(\mathbf{y}^{j})$$

$$t(\mathbf{y}^*)$$

(4)

(5)

$$\tilde{\mathbf{y}}(t + \Delta t) = \tilde{\mathbf{y}}(t) + \Delta t \frac{\mathrm{d}\mathbf{y}}{\mathrm{d}t}(t) + \frac{\Delta t^2}{2!} \frac{\mathrm{d}^2\mathbf{y}}{\mathrm{d}t^2}(t) + \cdots$$

$$\frac{1}{2}(t)+\cdots$$

$$\tilde{\mathbf{y}}(t + \Delta t) = \tilde{\mathbf{y}}(t) + \Delta t \, \tilde{\mathbf{f}}(\tilde{\mathbf{y}}(t)) + \frac{\Delta t^2}{2!} \left( \mathbf{D} \tilde{\mathbf{f}} \tilde{\mathbf{f}} \right) (\tilde{\mathbf{y}}(t)) + \cdots$$

### Example (cont.)

$$\tilde{\mathbf{y}}(t + \Delta t) = \tilde{\mathbf{y}}(t) + \Delta t \,\tilde{\mathbf{f}}(\tilde{\mathbf{y}}(t)) + \frac{\Delta t^2}{2!} \left(\mathbf{D}\tilde{\mathbf{f}}\tilde{\mathbf{f}}\right) (\tilde{\mathbf{y}}(t)) + \cdots$$

Substituting the ansatz:

$$\tilde{\mathbf{y}}(t + \Delta t) = \tilde{\mathbf{y}}(t) + \Delta t \, \mathbf{f} \left( \tilde{\mathbf{y}}(t) \right) + \Delta t^2 \left( \mathbf{f}_1 + \frac{1}{2} \mathbf{Dff} \right) \left( \tilde{\mathbf{y}}(t) \right) + \cdots$$
Set  $t = t^j$ , apply the defining condition  $\tilde{\mathbf{y}}(t^j) = \mathbf{y}^j$ ,  $\forall j$ ,

$$\mathbf{y}^{j+1} = \mathbf{y}^j + \Delta t \, \mathbf{f}(\mathbf{y}^j) + \Delta t^2 \left(\mathbf{f}_1 + \frac{1}{2} \mathbf{D} \mathbf{f} \mathbf{f}\right) \left(\mathbf{y}^j\right) + \cdots$$

and compare with the formula for the EE method:

$$\mathbf{y}^{j+1} = \mathbf{y}^j + \Delta t \, \mathbf{f}(\mathbf{y}^j)$$

$$\Rightarrow \mathbf{f}_1 = -rac{1}{2}\mathbf{Dff}$$
 ,  $\mathbf{f}_2 = \ldots$  , etc.

(6)

(7)

(8)

### Example (cont.)

DVF of the EE method in general (consistent, first-order):

$$\tilde{\mathbf{f}} = \mathbf{f} - \frac{\Delta t}{2} \mathbf{D} \mathbf{f} \mathbf{f} + \frac{\Delta t^2}{12} \mathbf{D}^2 (\mathbf{f}, \mathbf{f}) + \frac{\Delta t^2}{3} \mathbf{D} \mathbf{f} \mathbf{D} \mathbf{f} \mathbf{f} + \cdots$$
(9)

Mass-spring system (harmonic oscillator):

$$m\ddot{x}(t) + kx(t) = 0 \quad \Leftrightarrow \quad \underbrace{\begin{pmatrix} \dot{x} \\ \dot{v} \end{pmatrix}}_{\dot{y}} = \underbrace{\begin{pmatrix} 0 & 1 \\ -\omega^2 & 0 \end{pmatrix}}_{\dot{\mathbf{A}}} \underbrace{\begin{pmatrix} x \\ v \end{pmatrix}}_{\mathbf{y}}, \quad \text{with } \omega := \sqrt{\frac{k}{m}}$$
 (10)

DVF of this system simulated via the EE method:

$$\dot{\tilde{x}} = \tilde{v} + \frac{\Delta t}{2}\omega^2 \tilde{x} - \frac{\Delta t^2}{3}\omega^2 \tilde{v} + \mathcal{O}(\Delta t^3),$$

$$\dot{\tilde{v}} = -\omega^2 \tilde{x} + \frac{\Delta t}{2} \omega^2 \tilde{v} + \frac{\Delta t^2}{3} \omega^4 \tilde{x} + \mathcal{O}(\Delta t^3).$$

(12)

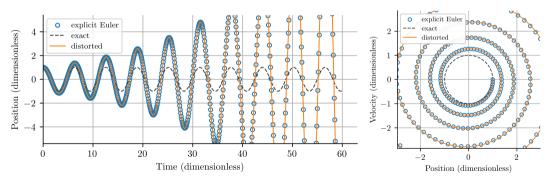
(11)

12/38

### Example (cont.)

Solution of DVF truncated up to  $\mathcal{O}(\Delta t^2)$ , position:

$$\tilde{x}(t) = C_1 e^{\frac{\Delta t}{2}\omega^2 t} \cos\left[\left(\omega - \frac{\Delta t^2}{3}\omega^3\right)t\right] + C_2 e^{\frac{\Delta t}{2}\omega^2 t} \sin\left[\left(\omega - \frac{\Delta t^2}{3}\omega^3\right)t\right], \quad (13)$$



→ notice: numerical anti-dissipation, shift in frequency (dispersion error)

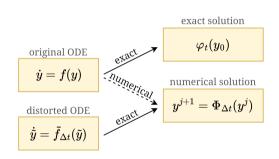
### BEA in practice: how to do it?

- Some more examples of this approach based on Taylor series: [3]
- ullet Equivalent, but different approaches also exist [4, 5, 6] o sometimes better suited for proofs
- Usually, the calculation quickly becomes tedious
- Quite algorithmic: original computer algebra implementation [3], extended in [7]

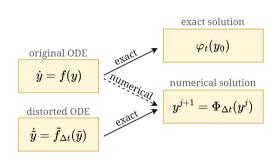
BEA-based compensation I. -

Modifying system parameters

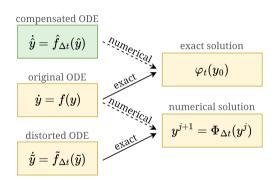
 BEA: tool for the analysis of structure-preserving numerical methods



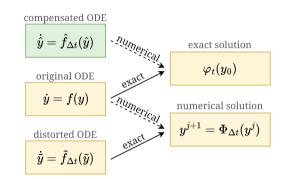
- BEA: tool for the analysis of structure-preserving numerical methods
- Can it also be applied constructively?



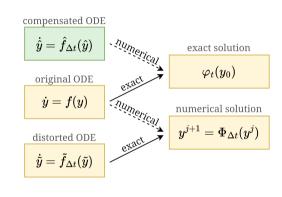
- BEA: tool for the analysis of structure-preserving numerical methods
- Can it also be applied constructively?
- Compensating the system being simulated



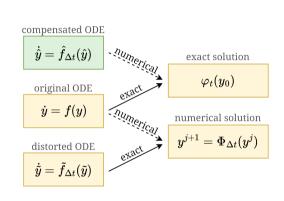
- BEA: tool for the analysis of structure-preserving numerical methods
- Can it also be applied constructively?
- Compensating the system being simulated
- Compensation purely through system parameters?



- BEA: tool for the analysis of structure-preserving numerical methods
- Can it also be applied constructively?
- Compensating the system being simulated
- Compensation purely through system parameters?
- Potential advantages:
  - Rectification of existing methods
  - No need to modify existing software implementations



- BEA: tool for the analysis of structure-preserving numerical methods
- Can it also be applied constructively?
- Compensating the system being simulated
- Compensation purely through system parameters?
- Potential advantages:
  - Rectification of existing methods
  - No need to modify existing software implementations
- Existing similar approaches [8, 9, 10] modify the vector field directly, not the parameters



- Demonstration: compensating the Newmark method [7]
- Newmark method: widely used time integration method for dynamics, available in commercial engineering software (ANSYS, Abaqus)
- Time-dependent finite element method (FEM) simulations for dynamics (direct time integration)
  - Wave propagation, crash simulations...
- Generally, numerical dissipation is present (absent only in special cases) – can be a (dis)advantage
- BEA-based compensation: opportunity to achieve better results with existing software
- ightarrow More accurate / less computationally intensive, more reliable numerical simulations





ODE: second-order linear system, corresponding to the semi-discrete equation of motion from the FEM model, with time-dependent external excitation

$$\mathbf{M\ddot{q}}(t) + \mathbf{C\dot{q}}(t) + \mathbf{Kq}(t) = \mathbf{F}(t), \quad \mathbf{q}(0) = \mathbf{q}_0, \, \dot{\mathbf{q}}(0) = \mathbf{v}_0. \tag{14}$$

To solve this numerically, the Newmark method is applied:

$$\mathbf{M}\mathbf{a}^{j+1} + \mathbf{C}\mathbf{v}^{j+1} + \mathbf{K}\mathbf{q}^{j+1} = \mathbf{F}^{j+1},$$
 (15)

$$\Delta t^2$$

$$\mathbf{q}^{j+1} = \mathbf{q}^j + \Delta t \mathbf{v}^j + \frac{\Delta t^2}{2} \left[ (1 - 2\beta) \mathbf{a}^j + 2\beta \mathbf{a}^{j+1} \right],$$

where 
$$\gamma,\,\beta$$
 are parameters of the Newmark method.  $\gamma=1/2\to {\rm symmetric},$ 

where  $\gamma, \beta$  are parameters of the Newmark method.  $\gamma = 1/2 \rightarrow$  symmetric, second-order

 $\mathbf{v}^{j+1} = \mathbf{v}^j + \Delta t \left[ (1 - \gamma) \mathbf{a}^j + \gamma \mathbf{a}^{j+1} \right].$ 

(16)

(17)

Rewriting (14) to an autonomous, first-order form is needed for BEA:

$$\underbrace{\begin{pmatrix} \dot{\mathbf{r}} \\ \dot{\mathbf{q}} \\ \dot{\mathbf{v}} \end{pmatrix}}_{\dot{\mathbf{y}}} = \underbrace{\begin{pmatrix} 1 \\ \mathbf{v} \\ -\mathbf{M}^{-1} \left( \mathbf{C} \mathbf{v} + \mathbf{K} \mathbf{q} - \mathbf{F}(\tau) \right) \end{pmatrix}}_{\mathbf{f}(\mathbf{y})}, \tag{18}$$

- $\Rightarrow$  the corresponding DVF  $ilde{\mathbf{f}}$  has been derived (margin on this slide too narrow)
- $\Rightarrow$  remarkably, the resulting DVF can be re-written in a second-order form with distorted matrices

$$\mathbf{M}\ddot{\mathbf{q}}(t) + \tilde{\mathbf{C}}\dot{\mathbf{q}}(t) + \tilde{\mathbf{K}}\mathbf{q}(t) = \tilde{\mathbf{F}}(t),$$
 (19)

⇒ this means that the Newmark method simulates a very similar system with different parameters: opportunity for compensation!

## Distorted matrices

$$\tilde{\mathbf{C}} = \mathbf{C} + \mathbf{W}(\Delta t, \gamma) \left( \mathbf{M}^{-1} \mathbf{K} - \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{C} \right) + \Delta t^{2} \left( \eta - \frac{1}{12} \right) \mathbf{K} \mathbf{M}^{-1} \mathbf{C}, 
\tilde{\mathbf{K}} = \mathbf{K} - \mathbf{W}(\Delta t, \gamma) \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{K} + \Delta t^{2} \left( \eta - \frac{1}{12} \right) \mathbf{K} \mathbf{M}^{-1} \mathbf{K}, 
\tilde{\mathbf{F}}(t) = \mathbf{F}(t) - \mathbf{W}(\Delta t, \gamma) \left( \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{F}(t) - \mathbf{M}^{-1} \mathbf{F}'(t) \right) + 
+ \Delta t^{2} \left( \eta - \frac{1}{12} \right) \left( \mathbf{K} \mathbf{M}^{-1} \mathbf{F}(t) - \mathbf{F}''(t) \right),$$

up to  $\mathcal{O}(\Delta t^2)$ , with

$$\eta = \frac{1}{2}\gamma - \beta - \frac{1}{12}, \quad \mathbf{W}(\Delta t, \gamma) = \Delta t \left(\gamma - \frac{1}{2}\right) \mathbf{M} - \Delta t^2 \left(\left(\gamma - \frac{1}{2}\right)^2 + \frac{1}{12}\right) \mathbf{C},$$
 and distorted initial conditions 
$$\mathbf{g}(0) = \mathbf{g}_0.$$

$$\mathbf{q}(0) = \mathbf{q}_0,$$

$$\dot{\mathbf{q}}(0) = \mathbf{q}_0 + \Delta t^2 \mathbf{q} \left( \mathbf{q} \right)$$

$$\mathbf{q}(0) = \mathbf{q}_0,$$

$$\dot{\mathbf{q}}(0) = \mathbf{v}_0 + \Delta t^2 \eta \left( -\mathbf{M}^{-1}\mathbf{C}\mathbf{M}^{-1}\mathbf{K}\mathbf{q}_0 + \left(\mathbf{M}^{-1}\mathbf{K} - \mathbf{M}^{-1}\mathbf{C}\mathbf{M}^{-1}\mathbf{C}\right)\mathbf{v}_0 + \right)$$

 $+ \mathbf{M}^{-1}\mathbf{C}\mathbf{M}^{-1}\mathbf{F}(0) - \mathbf{M}^{-1}\mathbf{F}'(0)$ .

$$\mathbf{K} - \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{C} \big) \, \mathbf{v}_0 \, +$$

(25)

(20)

(21)

(22)

(24)

19/38

Compensated system:

$$\mathbf{M}\ddot{\mathbf{q}}(t) + \hat{\mathbf{C}}\dot{\mathbf{q}}(t) + \hat{\mathbf{K}}\mathbf{q}(t) = \hat{\mathbf{F}}(t), \qquad \mathbf{q}(0) = \mathbf{q}_0, \quad \dot{\mathbf{q}}(0) = \mathbf{v}_0,$$
 (26)

Two compensations introduced:

- Eliminating numerical damping (not shown here)
- Achieving fourth-order accuracy
  - ullet for  $\gamma=1/2$ , only the  $\mathcal{O}\big(\Delta t^2\big)$  terms need to be eliminated

### Fourth-order compensation of the Newmark method

Result of derivation for the fourth-order compensation:

$$\widehat{\mathbf{C}} = \mathbf{C} + \frac{1}{12} \Delta t^2 \left( \mathbf{C} \mathbf{M}^{-1} \mathbf{K} + \mathbf{K} \mathbf{M}^{-1} \mathbf{C} - \mathbf{C} \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{C} \right), \tag{27}$$

$$\widehat{\mathbf{K}} = \mathbf{K} + \frac{1}{12} \Delta t^2 \left( \mathbf{K} \mathbf{M}^{-1} \mathbf{K} - \mathbf{C} \mathbf{M}^{-1} \mathbf{C} \mathbf{M}^{-1} \mathbf{K} \right),$$
(28)

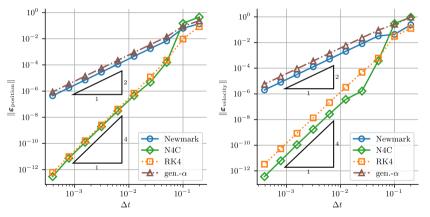
$$\widehat{\mathbf{F}}(t) = \mathbf{F}(t) + \frac{1}{12}\Delta t^2 \left( \mathbf{C}\mathbf{M}^{-1} \left( \mathbf{C}\mathbf{M}^{-1} \mathbf{F}(t) - \mathbf{F}'(t) \right) - \mathbf{K}\mathbf{M}^{-1} \mathbf{F}(t) + \mathbf{F}''(t) \right), \quad (29)$$

and  $\gamma = 1/2$ ,  $\beta = 1/6$  is required.

Once calculated, these can be used through the entire computation.

### Fourth-order compensation of the Newmark method

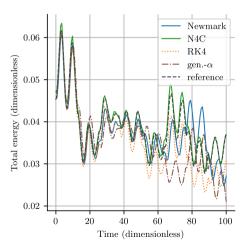
Convergence with  $\Delta t$  is indeed fourth-order:



Derivatives of  $\mathbf{F}(t)$  can also be calculated/estimated numerically with an at least second-order accurate formula.

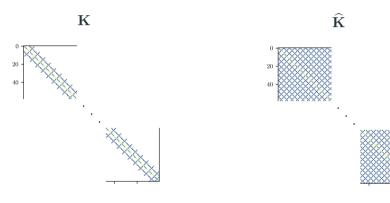
### Fourth-order compensation of the Newmark method

Non-continuous excitation (square wave), total energy of the system:



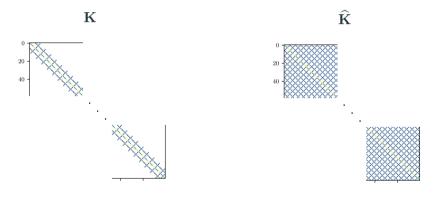
### Structure of the compensated matrices

 ${f M}$  does not change in the compensated system, but  ${f K}\to \widehat{{f K}}$  does – how does this change the sparsity structure of the FEM matrices?



### Structure of the compensated matrices

 ${f M}$  does not change in the compensated system, but  ${f K}\to \widehat{{f K}}$  does – how does this change the sparsity structure of the FEM matrices?

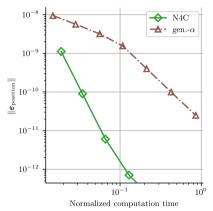


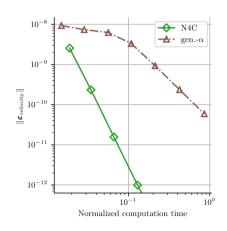
The derived DVF is underdetermined up to a matrix multiplier of  ${f M}$ 

 $\rightarrow$  opportunity for a better structure

### Structure of the compensated matrices

### Computation time vs. accuracy ( $\sim750~\text{DoF}$ )





### Open questions, future work

#### Open questions, future work

- How does the computational time / accuracy tradeoff scale with the number of DoFs? (ongoing work with D. Borza)
- ullet How do the eigenfrequencies change from  ${f K}$  to  $\widehat{{f K}}$ ?
- Can this approach be extended for more complex (nonlinear) excitation?
- Can a more favourable rewriting into second-order form be achieved with respect to the sparsity structure?
- $\bullet$  Extension of this approach to extensions of the Newmark method: HHT- $\alpha$  , generalised- $\alpha$

More details on this topic: [7] D. M. Takács and T. Fülöp. "Improving the accuracy of the Newmark method through backward error analysis". In: *Computational Mechanics* 75.5 (2025), pp. 1585–1606.

## Choosing a better coordinate

system

BEA-based compensation II. -

### Accuracy of a numerical method and the applied coordinate system

- We have seen so far that, through BEA, the DVF corresponding to a numerical method as applied to a system can be constructed
- For structure-preserving methods: the DVF is associated with a system of the same structure (detailed above)
- $\bullet$  Specifically, for symplectic numerical methods: the DVF is Hamiltonian; there is a distorted Hamiltonian  $\tilde{H}$  associated
- But: is the DVF invariant to a coordinate transformation of the original system?
  - Is the distorted Hamiltonian?
- Can this be exploited to achieve better accuracy?

# Hamiltonian systems, symplectic methods, coordinate transformations

Hamiltonian system (autonomous): described by the Hamiltonian  $H(\mathbf{q}, \mathbf{p})$ ,  $H: \mathbb{R}^d \times \mathbb{R}^d \to \mathbb{R}$ 

Equations of motion:

$$\underbrace{\begin{pmatrix} \dot{\mathbf{q}} \\ \dot{\mathbf{p}} \end{pmatrix}}_{\dot{\mathbf{y}}} = \underbrace{\begin{pmatrix} \mathbf{0} & \mathbf{I} \\ -\mathbf{I} & \mathbf{0} \end{pmatrix}}_{\mathcal{J}_{c}} \underbrace{\begin{pmatrix} H_{q}(\mathbf{q}, \mathbf{p}) \\ H_{p}(\mathbf{q}, \mathbf{p}) \end{pmatrix}}_{\mathbf{D}H(\mathbf{y})},$$

Symplectic property:

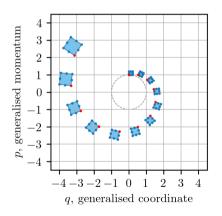
od: 
$$\left( \frac{\partial \Phi_{\Delta t}}{\partial \mathbf{v}} \right)^\mathsf{T} \mathcal{J}_\mathrm{c}^{-1} \left( \frac{\partial \Phi_{\Delta t}}{\partial \mathbf{v}} \right) = \mathcal{J}_\mathrm{c}^{-1}$$

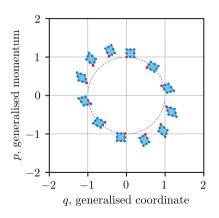
$$\left(\frac{\partial \varphi_t}{\partial \mathbf{v}}\right)^{\mathsf{T}} \mathcal{J}_{\mathbf{c}}^{-1} \left(\frac{\partial \varphi_t}{\partial \mathbf{v}}\right) = \mathcal{J}_{\mathbf{c}}^{-1} \tag{33}$$

(30)

### Hamiltonian systems, symplectic methods, coordinate transformations

Explicit Euler vs. Symplectic Euler method, harmonic oscillator:





### Hamiltonian systems, symplectic methods, coordinate transformations

To canonical transformations

$$\begin{pmatrix} \mathbf{q} \\ \mathbf{p} \end{pmatrix} \mapsto \begin{pmatrix} \bar{\mathbf{q}} \\ \bar{\mathbf{p}} \end{pmatrix} = \begin{pmatrix} \mathbf{Q}(\mathbf{q}, \mathbf{p}) \\ \mathbf{P}(\mathbf{q}, \mathbf{p}) \end{pmatrix}, \tag{33}$$

the Hamiltonian is invariant, i.e.

$$\bar{H}(\bar{\mathbf{q}}, \bar{\mathbf{p}}) := H(\mathbf{Q}^{-1}(\bar{\mathbf{q}}, \bar{\mathbf{p}}), \mathbf{P}^{-1}(\bar{\mathbf{q}}, \bar{\mathbf{p}})); \qquad \bar{H} = H$$
 (34)

From here on, we restrict ourselves to the subset of canonical transformations induced by a coordinate transformation  $\bar{\mathbf{q}}=\mathbf{Q}(\mathbf{q})$ :

$$\bar{q}^{\alpha} = Q^{\alpha}(\mathbf{q}); \quad \bar{p}_{\alpha} = p_i \left( \frac{\partial \left( Q^{-1} \right)^i}{\partial \bar{q}^{\alpha}} \right) \circ \mathbf{Q}(\mathbf{q})$$
 (35)

Does this change of coordinates change the numerical results in a meaningful way?

### Distorted Hamiltonian of a symplectic method

Symplectic Euler (SE) method (first-order):

$$\mathbf{q}^{j+1} = \mathbf{q}^j + \Delta t H_p(\mathbf{q}^{j+1}, \mathbf{p}^j),$$
  
 $\mathbf{p}^{j+1} = \mathbf{p}^j - \Delta t H_q(\mathbf{q}^{j+1}, \mathbf{p}^j).$ 

Distorted Hamiltonian from BEA:

$$\tilde{H} = H - \frac{\Delta t}{2} H_p H_q + \mathcal{O}(\Delta t^2),$$

Similarly, in the transformed coordinate system:

~ 
$$\Delta t$$

 $\tilde{\bar{H}} = \bar{H} - \frac{\Delta t}{2} \bar{H}_{\bar{p}} \bar{H}_{\bar{q}} + \mathcal{O}(\Delta t^2).$ 

$$\tilde{\bar{H}} = \bar{H} - \frac{\Delta t}{4} \bar{H}_{\bar{z}} \bar{H}_{\bar{z}}$$

 $\rightarrow$  are these actually the same, i.e., is H invariant to coordinate transformations?

(38)

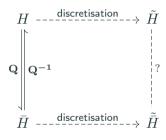
(36)

(37)

$$ilde{ar{H}} = ar{H} - rac{\Delta t}{L} ar{H}_{ar{n}} ar{H}_{ar{a}}$$

31/38

### Distorted Hamiltonian of a symplectic method



### Distorted Hamiltonian of a symplectic method

$$\tilde{\bar{H}} = \bar{H} - \frac{\Delta t}{2} \bar{H}_{\bar{p}} \bar{H}_{\bar{q}} + \mathcal{O}(\Delta t^2). \tag{38}$$

Let us consider the second term in (38),  $H_pH_q$  (first elementary Hamiltonian):

$$\bar{H}_{\bar{p}}\bar{H}_{\bar{q}} \equiv \frac{\partial \bar{H}}{\partial \bar{p}_{\alpha}} \frac{\partial \bar{H}}{\partial \bar{q}^{\alpha}} = \dots = \underbrace{\frac{\partial H}{\partial p_{i}} \frac{\partial H}{\partial p_{k}} p_{l} \frac{\partial \left(Q^{-1}\right)^{l}}{\partial \bar{q}^{\beta}} \frac{\partial^{2} Q^{\beta}}{\partial q^{k} \partial q^{i}}}_{=:\Xi_{H_{p}H_{q},Q}} + \frac{\partial H}{\partial p_{i}} \frac{\partial H}{\partial q^{i}}.$$
(39)

→ necessary condition for the invariance of the distorted Hamiltonian:

$$\Xi_{H_pH_q,\mathbf{Q}}=0$$
  $\bullet$  trivially fulfilled for affine coordinate transformations, can also be fulfilled

non-invariant in general

nontrivially

- can this be exploited to eliminate the first-order term in H to raise accuracy?

(40)

### Second-order accuracy with the symplectic Euler method

Condition for second-order accuracy of SE in the transformed coordinate system:

$$\bar{H}_{\bar{p}}\bar{H}_{\bar{q}} \equiv 0 \quad \Leftrightarrow \quad \Xi_{H_pH_q,\mathbf{Q}} = -\frac{\partial H}{\partial p_i}\frac{\partial H}{\partial q^i}$$
 (41)

- trivially fulfilled in fully cyclic (action-angle) coordinates.
  - achieving this is usually not possible
- ullet are there any other, non-trivial  ${f Q}$  coordinate transformations that achieve this?
  - not guaranteed, but sometimes possible

### Demonstration: harmonic oscillator, symplectic Euler method

Harmonic oscillator:

$$H(q,p) = \frac{1}{2}p^2 + \frac{1}{2}q^2,$$

Condition (41) for second-order SE in this case, after calculations:

$$(1-q^2)\left(\frac{\partial Q}{\partial q}\right)^{-1}\frac{\partial^2 Q}{\partial q^2} + q = 0.$$

Appropriate coordinate transformation fulfilling this condition:

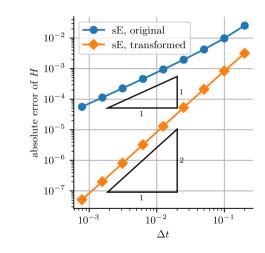
$$Q(q) = \frac{2}{\pi} \left( \hat{q} \sqrt{1 - \hat{q}^2} + \arcsin(\hat{q}) \right), \qquad \text{where } \hat{q} = \frac{q}{1 + 2\Delta t^2},$$

(42)

(43)

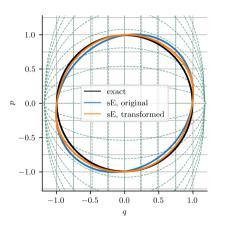
(44)

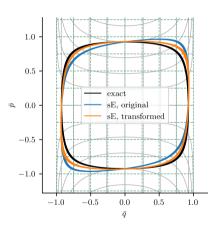
### Demonstration



#### Demonstration

Trajectory of the numerical results in the two coordinate systems:





### Open questions, future work, other results

- When does a better coordinate system exist in general?
- Extension of this approach beyond the SE method? (Størmer–Verlet might be a good candidate)
- Extension to the entire class of canonical transformations?
- Additional result not explored here: preservation of first integrals in the SE method

More details on this topic: [11] D. M. Takács and T. Fülöp. "On the coordinate system-dependence of the accuracy of symplectic methods". In: *Journal of Numerical Analysis and Approximation Theory* (2025). In press.

### Bibliography i

- [1] P. C. Moan. *On rigorous modified equations for discretizations of ODEs.*Tech. rep. Technical Report 2005-3, Geometric Integration Preprint Server, 2005.
- [2] G. Benettin and A. Giorgilli. "On the Hamiltonian interpolation of near-to-the identity symplectic mappings with application to symplectic integration algorithms". In: *Journal of Statistical Physics* 74 (1994), pp. 1117–1143.
- [3] E. Hairer and C. Lubich. "Asymptotic Expansions and Backward Analysis for Numerical Integrators". In: *Dynamics of Algorithms*. Springer New York, 2000, pp. 91–106.
- [4] S. Reich. "Backward error analysis for numerical integrators". In: SIAM Journal on Numerical Analysis 36.5 (1999), pp. 1549–1570.

### Bibliography ii

- O. Gonzalez, D. J. Higham, and A. M. Stuart. "Qualitative properties of modified equations". In: *IMA Journal of Numerical Analysis* 19.2 (1999), pp. 169–190.
- [6] P. C. Moan. "On modified equations for discretizations of ODEs". In: Journal of Physics A: Mathematical and General 39.19 (2006), p. 5545.
- [7] D. M. Takács and T. Fülöp. "Improving the accuracy of the Newmark method through backward error analysis". In: Computational Mechanics 75.5 (2025), pp. 1585–1606.
- [8] B. A. Shadwick, J. C. Bowman, and P. J. Morrison. "Exactly conservative integrators". In: SIAM Journal on Applied Mathematics 59.3 (1998), pp. 1112–1133.

### Bibliography iii

- [9] X. Shang and H. C. Öttinger. "Structure-preserving integrators for dissipative systems based on reversible-irreversible splitting". In: Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences 476.2234 (Feb. 2020), p. 20190446.
- [10] P. Chartier, E. Hairer, and G. Vilmart. "Numerical integrators based on modified differential equations". In: Mathematics of computation 76.260 (2007), pp. 1941–1953.
- [11] D. M. Takács and T. Fülöp. "On the coordinate system-dependence of the accuracy of symplectic methods". In: Journal of Numerical Analysis and Approximation Theory (2025). In press.

Thank you for your kind attention!

### Non-invariance of the distorted Hamiltonian

