ADER and DeC:
arbitrarily high order (explicit)
methods for PDEs and ODEs

Davide Torlo

*MathLab, Mathematics Area, SISSA International
School for Advanced Studies, Trieste, Italy
davidetorlo.it

## Outline

## Outline

ARBITURLY

<u>Motivation: high order accurate (explicit) method</u>

Methods used to solve a hyperbolic PDE system for $u : \mathbb{R}^+ \times \Omega \to \mathbb{R}^D$

$$\partial_t u + \nabla_x \mathcal{F}(u) = 0. \tag{1}$$

Or ODE system for $\boldsymbol{u} : \mathbb{R}^+ \to \mathbb{R}^S$

$$\partial_t \boldsymbol{u} = F(\boldsymbol{u}). \tag{2}$$

Applications:

- Fluids/transport
- Chemical/biological processes

How?

- <u>Arbitrarily</u> <u>high order</u> accurate
-

Methods used to solve a hyperbolic PDE system for $u : \mathbb{R}^+ \times \Omega \to \mathbb{R}^D$

(1)

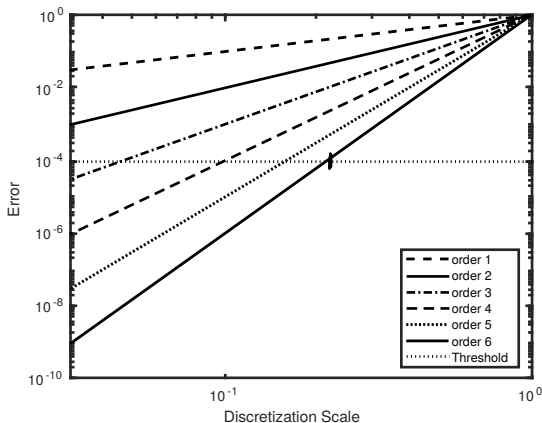Or ODE system for $\boldsymbol{u}$ :

(2)

Applications:

- Fluids/transport
- Chemical/biological

How?

- Arbitrarily high order
-

Methods used to solve a hyperbolic PDE system for $u : \mathbb{R}^+ \times \Omega \to \mathbb{R}^D$

$$\partial_t u + \nabla_{\mathbf{x}} \mathcal{F}(u) = 0. \tag{1}$$

Or ODE system for $\boldsymbol{u} : \mathbb{R}^+ \to \mathbb{R}^S$

$$\partial_t \boldsymbol{u} = F(\boldsymbol{u}). \tag{2}$$

Applications:
- Fluids/transport
- Chemical/biological processes

How?
- Arbitrarily high order accurate
- Explicit (if nonstiff problem)

## Classical time integration: Runge–Kutta

$$
\begin{cases}
\boldsymbol{u}^{(1)} := \boldsymbol{u}^n, & (3) \\[2ex]
\boldsymbol{u}^{(k)} := \boldsymbol{u}^n + \sum_{s=1}^{K} a_{ks} F\left(t^n + c_s \Delta t, \boldsymbol{u}^{(s)}\right), \quad \text{for } k = 2, \ldots, K, & (4) \\[2ex]
\boldsymbol{u}^{n+1} := \boldsymbol{u}^n + \sum_{s=1}^{K} b_s F\left(t^n + c_s \Delta t, \boldsymbol{u}^{(s)}\right). & (5)
\end{cases}
$$

$$\boldsymbol{u}^{(k)} := \boldsymbol{u}^n + \sum_{s=1}^{k-1} a_{ks} F\left(t^n + c_s \Delta t, \boldsymbol{u}^{(s)}\right), \quad \text{for } k = 2, \ldots, K.$$

- Easy to solve
- High orders involved:
  - Order conditions: system of many equations
  - Stages $K \geq d$ order of accuracy (e.g. RK44, RK65)

$$\boldsymbol{u}^{(k)} := \boldsymbol{u}^n + \sum_{s=1}^{K} a_{ks} F\left(t^n + c_s \Delta t, \boldsymbol{u}^{(s)}\right), \quad \text{for } k = 2, \dots, K.$$

- More complicated to solve for nonlinear systems
- High orders easily done:
  - Take a high order quadrature rule on $[t^n, t^{n+1}]$
  - Compute the coefficients accordingly, see Gauss–Legendre or Gauss–Lobatto polynomials
  - Order up to $d = 2K$

## ADER and DeC

Two iterative explicit arbitrarily high order accurate methods.

- <u>ADER</u>[1] for hy<u>perbolic PDE</u>, after a first analytic more complicated approach.
- Deferred Correction (DeC): introduced for <u>explicit ODE</u>[2], extended to implicit ODE[3] and to hyperbolic PDE[4].

---

[1]M. Dumbser, D. S. Balsara, E. F. Toro, and C.-D. Munz. A unified framework for the construction of one-step finite volume and discontinuous galerkin schemes on unstructured meshes. Journal of Computational Physics, 227(18):8209–8253, 2008.

[2]A. Dutt, L. Greengard, and V. Rokhlin. Spectral Deferred Correction Methods for Ordinary Differential Equations. BIT Numerical Mathematics, 40(2):241–266, 2000.

[3]M. L. Minion. Semi-implicit spectral deferred correction methods for ordinary differential equations. Commun. Math. Sci., 1(3):471–500, 09 2003.

[4]R. Abgrall. High order schemes for hyperbolic problems using globally continuous approximation and avoiding mass matrices. Journal of Scientific Computing, 73(2):461–494, Dec 2017.

# Outline

High order in time: we discretize our variable on $[t^n, t^{n+1}]$ in $M$ substeps ($u^m$).

$$\partial_t u = F(u(t)).$$

Thanks to Picard–Lindelöf theorem, we can rewrite

INTEGRAL
FORM

$$u^m = u^0 + \int_{t^0}^{t^m} F(u(t))dt.$$

and if we want to reach order $r+1$ we need $M = r$.

- EQUIDISTANT POINTS
- GAUSS - LOBATTO

$\varphi_n$ LAGRANGIAN POLYS

$u^m = u^0 + \int_{t^0}^{t^n} F(u(t)) \, dt$

$\varphi_n(t^n) = \delta_{m n}$

$u^M \quad t^M$

More precisely, ~~we want to~~ we want to solve $\mathcal{L}^2(u^{n,0}, \ldots, u^{n,M}) = 0$, where

$$\mathcal{L}^2(u^0, \ldots, u^M) = \begin{pmatrix} u^M - u^0 + \sum_{r=0}^{M} \int_{t^0}^{t^M} F(u^r)\varphi_r(s)\mathrm{d}s \\ \vdots \\ u^1 - u^0 + \sum_{r=0}^{M} \int_{t^0}^{t^1} F(u^r)\varphi_r(s)\mathrm{d}s \end{pmatrix}$$

$\Theta_n^m = \int_0^{t^n} \varphi_n(s) \, ds$

$u^m \quad t^m$

$v \in \mathbb{R}^S$

- $\mathcal{L}^2 = 0$ is a system of $\boxed{M} \times \boxed{S}$ coupled (non)linear equations
- $\mathcal{L}^2$ is an implicit method (collocation method: Gauss, LobattoIIIA)
- Not easy to solve directly $\mathcal{L}^2(\underline{u}^*) = 0$
- High order ($\geq M + 1$), depending on points distribution

$u^1 \quad t^1$

$u^0 \quad t^0$

More precisely, for each $\sigma$ we want to solve $\mathcal{L}^2(\boldsymbol{u}^{n,0}, \ldots, \boldsymbol{u}^{n,M}) = 0$, where

$$\mathcal{L}^2(\boldsymbol{u}^0, \ldots, \boldsymbol{u}^M) = \begin{pmatrix} \boldsymbol{u}^M - \boldsymbol{u}^0 + \Delta t \sum_{r=0}^M \theta_r^M F(\boldsymbol{u}^r) \\ \vdots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 + \Delta t \sum_{r=0}^M \theta_r^1 F(\boldsymbol{u}^r) \end{pmatrix} = 0$$

$$\Theta_n = \underbrace{\frac{1}{\Delta t} \int_{t^0}^{t^m} \varphi_n(s) \, ds}_{} = \int_0^{\frac{t^m}{t^n}} \widetilde{\varphi}_n(s) \, ds$$

$$U^m = U^0 + \Delta t \sum \Theta_n^m F(U^n) \quad \forall m$$

- $\mathcal{L}^2 = 0$ is a system of $\underline{M} \times \underline{S}$ coupled (non)linear equations
- $\mathcal{L}^2$ is an implicit method (collocation method: Gauss, LobattoIIIA)
- Not easy to solve directly $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$ $\quad u^* \text{ UNKNOWN}$
- High order ($\geq \underline{M+1}$), depending on points distribution

$$2M \quad \text{LOBATTO}$$

$$\partial_t u = F(u) \qquad \partial_t u + F(u) = 0 \qquad u^M \quad t^M$$

Instead of solving the implicit system directly (difficult), we introduce a first order scheme $\mathcal{L}^1(u^{n,0}, \ldots, u^{n,M})$:

EXPLICIT EULER

$$\mathcal{L}^1(u^0, \ldots, u^M) = \begin{pmatrix} u^M - u^0 + \Delta t \beta^M F(u^0) \\ \vdots \\ u^1 - u^0 + \Delta t \beta^1 F(u^0) \end{pmatrix}$$

- First order approximation
- Explicit Euler
- Easy to solve $\underline{\mathcal{L}^1(\underline{u}) = 0}$

$$\mathcal{L}'(\underline{u}) = \mathbb{1}$$

$$\beta^m \Delta t = t^m - t^0$$
$$\beta^m = \frac{t^m - t^0}{\Delta t}$$

$u^m \quad t^m \qquad \beta^m$

$u^1 \quad t^1$

$u^0 \quad t^0$

## Deferred Correction[5]

How to combine two methods keeping the accuracy of the second and the stability and simplicity of the first one?

$$\boxed{\underline{\boldsymbol{u}}^{0,(k)} := \boldsymbol{u}(t^n),} \quad k = 0, \ldots, K,$$
$$\boxed{\underline{\boldsymbol{u}}^{m,(0)} := \boldsymbol{u}(t^n),} \quad m = 1, \ldots, M$$

$\rightsquigarrow \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)})$ with $k = 1, \ldots, K$.

### Theorem (Convergence DeC)

- $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$
- If $\mathcal{L}^1$ coercive with constant $C_1$
- If $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz with constant $C_2 \Delta t$

Then $\|\underline{\boldsymbol{u}}^{(K)} - \underline{\boldsymbol{u}}^*\| \leq C(\Delta t)^K$

- $\mathcal{L}^1(\underline{\boldsymbol{u}}) = 0$, first order accuracy, easily invertible.
- $\mathcal{L}^2(\underline{\boldsymbol{u}}) = 0$, high order $M + 1$.



---

[5]A. Dutt, L. Greengard, and V. Rokhlin. BIT Numerical Mathematics, 40(2):241–266, 2000.

## Deferred Correction[5]

How to combine two methods keeping the accuracy of the second and the stability and simplicity of the first one?

- $\mathcal{L}^1(\underline{\boldsymbol{u}}) = 0$, first order accuracy, easily invertible.
- $\mathcal{L}^2(\underline{\boldsymbol{u}}) = 0$, high order $M + 1$.

$$\boldsymbol{u}^{0,(k)} := \boldsymbol{u}(t^n), \quad k = 0, \ldots, K,$$

$$\boldsymbol{u}^{m,(0)} := \boldsymbol{u}(t^n), \quad m = 1, \ldots, M$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}) \text{ with } k = 1, \ldots, K.$$

### Theorem (Convergence DeC)

- $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$
- If $\mathcal{L}^1$ coercive with constant $C_1$
- If $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz with constant $C_2 \Delta t$

Then $\|\underline{\boldsymbol{u}}^{(K)} - \underline{\boldsymbol{u}}^*\| \leq C(\Delta t)^K$



---

[5]A. Dutt, L. Greengard, and V. Rokhlin. BIT Numerical Mathematics, 40(2):241–266, 2000.

## Deferred Correction[5]

How to combine two methods keeping the accuracy of the second and the stability and simplicity of the first one?

- $\mathcal{L}^1(\underline{\boldsymbol{u}}) = 0$, first order accuracy, easily invertible.
- $\mathcal{L}^2(\underline{\boldsymbol{u}}) = 0$, high order $M + 1$.

$$\boldsymbol{u}^{0,(k)} := \boldsymbol{u}(t^n), \quad k = 0, \ldots, K,$$

$$\boldsymbol{u}^{m,(0)} := \boldsymbol{u}(t^n), \quad m = 1, \ldots, M$$

$\longrightarrow$ $\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)})$ with $k = 1, \ldots, K$.

### Theorem (Convergence DeC)

- $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$
- If $\mathcal{L}^1$ coercive with constant $C_1$
- If $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz with constant $C_2 \Delta t$

Then $\|\underline{\boldsymbol{u}}^{(K)} - \underline{\boldsymbol{u}}^*\| \leq C(\Delta t)^K$



---

[5]A. Dutt, L. Greengard, and V. Rokhlin. BIT Numerical Mathematics, 40(2):241–266, 2000.

## Deferred Correction[5]

How to combine two methods keeping the accuracy of the second and the stability and simplicity of the first one?

$$\boldsymbol{u}^{0,(k)} := \boldsymbol{u}(t^n), \quad k = 0, \ldots, K,$$

$$\boldsymbol{u}^{m,(0)} := \boldsymbol{u}(t^n), \quad m = 1, \ldots, M$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}) \text{ with } k = 1, \ldots, K.$$

- $\mathcal{L}^1(\underline{\boldsymbol{u}}) = 0$, first order accuracy, easily invertible.
- $\mathcal{L}^2(\underline{\boldsymbol{u}}) = 0$ high order $M + 1$.

### Theorem (Convergence DeC)

- $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$
- If $\mathcal{L}^1$ coercive with constant $C_1$
- If $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz with constant $C_2 \Delta t$

Then $\|\underline{\boldsymbol{u}}^{(K)} - \underline{\boldsymbol{u}}^*\| \leq C(\Delta t)^K$



[5]A. Dutt, L. Greengard, and V. Rokhlin. BIT Numerical Mathematics, 40(2):241–266, 2000.

## Deferred Correction[5]

How to combine two methods keeping the accuracy of the second and the stability and simplicity of the first one?

- $\mathcal{L}^1(\underline{\boldsymbol{u}}) = 0$, first order accuracy, easily invertible.
- $\mathcal{L}^2(\underline{\boldsymbol{u}}) = 0$, high order $M + 1$.

$$\boldsymbol{u}^{0,(k)} := \boldsymbol{u}(t^n), \quad k = 0, \dots, K,$$
$$\boldsymbol{u}^{m,(0)} := \boldsymbol{u}(t^n), \quad m = 1, \dots, M$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}) \text{ with } k = 1, \dots, K.$$

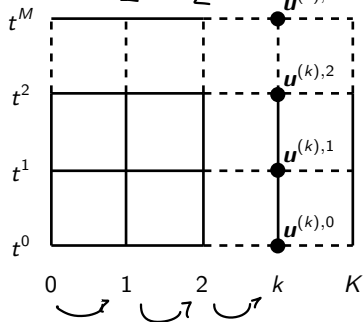### Theorem (Convergence DeC)

- $\mathcal{L}^2(\underline{\boldsymbol{u}}^*) = 0$
- If $\mathcal{L}^1$ coercive with constant $C_1$
- If $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz with constant $C_2\Delta t$

Then $\|\underline{\boldsymbol{u}}^{(K)} - \underline{\boldsymbol{u}}^*\| \leq C(\Delta t)^K$



---

[5]A. Dutt, L. Greengard, and V. Rokhlin. BIT Numerical Mathematics, 40(2):241–266, 2000.

THESIS:• $\| u^{(K)} - U^* \| \leq (c \Delta t)^K \cdot \| U^{(0)} - U^* \|$

HYP:in $\exists u^*; \mathcal{L}^2(u^*) = 0$    2) $\| \mathcal{L}^1(u) - \mathcal{L}^1(v) \| \geq c_1 \| u - v \|$    3) $\| \mathcal{L}^2(u) - \mathcal{L}^2(u) - (\mathcal{L}^2(v)$
$- \mathcal{L}^2(v)) \| \leq \| u - v \| \cdot c_2 \Delta t$

## Proof.

Let $\underline{u}^*$ be the solution of $\mathcal{L}^2(\underline{u}^*) = 0$. We know that $\mathcal{L}^1(\underline{u}^*) = \mathcal{L}^1(\underline{u}^*) - \mathcal{L}^2(\underline{u}^*)$, so that

$$\| U^{(K)} - U^* \| \leq \frac{1}{c_1} \| \mathcal{L}^1(U^{(K)}) - \mathcal{L}^1(U^*) \| = \frac{1}{c_1} \| \mathcal{L}^1(U^{(K-1)}) - \mathcal{L}^2(U^{(K-1)}) - \mathcal{L}^1(U^*) + \mathcal{L}^2(U^*) \|$$

$$\underset{c1P.}{\leq} \frac{c_2 \cdot \Delta t}{\underbrace{c_1}_{c \Delta t}} \| U^{(K-1)} - U^* \| \leq (c \Delta t)^2 \| U^{(K-2)} - U^* \| \leq \ldots \leq (c \Delta t)^K \| U^{(0)} - U^* \|$$

  □

ITERATIVE PROCESS

$$\mathcal{L}^1(U^{(K)}) = \mathcal{L}^1(U^{(K-1)}) - \mathcal{L}^2(U^{(K-1)})$$

**Proof.**

Let $f^*$ be the solution of $\mathcal{L}^2(\underline{u}^*) = 0$. We know that $\mathcal{L}^1(\underline{u}^*) = \mathcal{L}^1(\underline{u}^*) - \mathcal{L}^2(\underline{u}^*)$, so that

$$\mathcal{L}^1(\underline{u}^{(k+1)}) - \mathcal{L}^1(\underline{u}^*) = \left(\mathcal{L}^1(\underline{u}^{(k)}) - \mathcal{L}^2(\underline{u}^{(k)})\right) - \left(\mathcal{L}^1(\underline{u}^*) - \mathcal{L}^2(\underline{u}^*)\right)$$

$$C_1||\underline{u}^{(k+1)} - \underline{u}^*|| \leq ||\mathcal{L}^1(\underline{u}^{(k+1)}) - \mathcal{L}^1(\underline{u}^*)|| =$$

$$= ||\mathcal{L}^1(\underline{u}^{(k)}) - \mathcal{L}^2(\underline{u}^{(k)}) - (\mathcal{L}^1(\underline{u}^*) - \mathcal{L}^2(\underline{u}^*))|| \leq$$

$$\leq C_2\Delta||\underline{u}^{(k)} - \underline{u}^*||.$$

$$||\underline{u}^{(k+1)} - \underline{u}^*|| \leq \left(\frac{C_2}{C_1}\Delta\right)||\underline{u}^{(k)} - \underline{u}^*|| \leq \left(\frac{C_2}{C_1}\Delta\right)^{k+1}||\underline{u}^{(0)} - \underline{u}^*||.$$

After $K$ iteration we have an error at most of $\left(\frac{C_2}{C_1}\Delta\right)^K||\underline{u}^{(0)} - \underline{u}^*||$. $\quad\square$

COERCIVITY: $\| \mathcal{L}^2(u) - \mathcal{L}^2(v) \| \geq C_1 \| u - v \|$

$$\mathcal{L}^1(u) - \mathcal{L}^1(v) = \begin{pmatrix} u^n - u^o + \beta^n \Delta t F(v^o) \\ \vdots \\ u^2 - u^o + \beta^1 \Delta t F(u^o) \end{pmatrix} - \begin{pmatrix} v^n - u^o + \beta^n \Delta t F(u^o) \\ v^2 - u^o + \beta^1 \Delta t F(u^o) \end{pmatrix}$$

$$= \begin{pmatrix} u^n - v^n \\ u^1 - v^1 \end{pmatrix} = \underline{u} - \underline{v} \qquad C_1 = 1$$

LIPSCHITZ CONT

$\mathcal{L}^1$ EXPL EUL APPROX          $\mathcal{L}^2$ H.O. IMPLICIT RK METH

$\mathcal{L}^1(u^{ex}) = \underline{\mathcal{O}(\Delta t^2)}$          $\mathcal{L}^2(u^{ex}) = \mathcal{O}(\Delta t^{P+1})$

$|\mathcal{L}^1 - \mathcal{L}^2| = \mathcal{O}(\Delta t^2)$

SKETCH
PROOF WITH
DETAILS
EFFICIENT DeC
M. CALIZZI, TORLO

$t^1$

$\mathcal{L}^2(u) = \left( u^2 - u^0 + dt \frac{1}{2} \left( F(v^0) - F(v^1) \right) \right)$

$t^0$

$\mathcal{L}^2(u) = u^1 - u^0 + dt \, F(v^0)$

ITERATIVE PROCESS

$U^{(0),0} = U^{(0),1} = u(t^n) = U^0 \quad \boxed{k=0}$

$k=1 \quad \mathcal{L}^1(U^{(1)}) = \mathcal{L}^1(U^{(0)}) - \mathcal{L}^2(u^{(0)})$

$U^{(1),1} - \cancel{U^0 + dt \, F(v^0)} = \cancel{U^0 - U^0 + dt \, F(v^0)} - \left( \cancel{u^0} - U^0 + \frac{dt}{2} \left[ F(v^0) + F(v^0) \right] \right)$

$U^{(1),1} = U^0 - \Delta t \, F(U^0) \quad (\text{EXPLICIT EULER 1ST ORDER})$

$\underline{k=2} \quad \mathcal{L}^1(U^{(2)}) = \mathcal{L}^1(U^{(1)}) - \mathcal{L}^2(U^{(1)})$

$U^{(2),1} - \cancel{U^0 + dt \, F(v^0)} = \cancel{U^{(1),1} - U^0 + dt \, F(v^0)} - \left[ U^{(1),1} - U^0 + \frac{dt}{2} \left( F(v^0) + F(U^{(1),1}) \right) \right]$

$$U^{(2),1} = U^0 - \Delta \frac{t}{2} \left( \overline{F}(U^0) + \overline{F}(U^{(1),1}) \right)$$

$2^{nd}$ ORDER ACCURATE

# DeC: Second order example

In practice

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

For $m = 1, \ldots, M$

$$\boldsymbol{u}^{(k),m} \underline{\phantom{-}} \boldsymbol{u}^0 - \beta^m \Delta t F(\boldsymbol{u}^0) - \boldsymbol{u}^{(k-1),m} + \boldsymbol{u}^0 + \beta^m \Delta t F(\boldsymbol{u}^0)$$

$$+ \boldsymbol{u}^{(k-1),m} \underline{\phantom{-}} \boldsymbol{u}^0 - \Delta t \sum_{r=0}^{M} \theta_r^m F(\boldsymbol{u}^{(k-1),r}) = 0$$

In practice

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

For $m = 1, \ldots, M$

$$\boldsymbol{u}^{(k),m} \underline{-\, \boldsymbol{u}^0 - \beta^m \Delta t F(\boldsymbol{u}^0)} - \boldsymbol{u}^{(k-1),m} + \underline{\boldsymbol{u}^0 + \beta^m \Delta t F(\boldsymbol{u}^0)}$$

$$+ \boldsymbol{u}^{(k-1),m} \underline{-\, \boldsymbol{u}^0} - \Delta t \sum_{r=0}^{M} \theta_r^m F(\boldsymbol{u}^{(k-1),r}) = 0$$

In practice

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

For $m = 1, \ldots, M$

$$\boldsymbol{u}^{(k),m} \underline{\ \ } \underline{\boldsymbol{u}^0} \underline{\ -\ \beta^m \Delta t F(\boldsymbol{u}^0)} - \cancel{\boldsymbol{u}^{(k-1),m}} + \underline{\boldsymbol{u}^0 + \beta^m \Delta t F(\boldsymbol{u}^0)}$$

$$+ \cancel{\boldsymbol{u}^{(k-1),m}} \underline{\ \ } \boldsymbol{u}^0 - \Delta t \sum_{r=0}^{M} \theta_r^m F(\boldsymbol{u}^{(k-1),r}) = 0$$

In practice

$$\mathcal{L}^1(\underline{u}^{(k)}) = \mathcal{L}^1(\underline{u}^{(k-1)}) - \mathcal{L}^2(\underline{u}^{(k-1)}), \qquad k = 1, \ldots, K,$$

For $m = 1, \ldots, M$

$$u^{(k),m} \underbrace{u^0 - \beta^m \Delta t F(u^0)}_{M} - u^{(k-1),m} + u^0 + \beta^m \Delta t F(u^0)$$

$$+ u^{(k-1),m} u^0 - \Delta t \sum_{r=0}^{M} \theta_r^m F(u^{(k-1),r}) = 0$$

$$\boxed{u^{(k),m} - u^0 - \Delta t \sum_{r=0}^{M} \theta_r^m F(u^{(k-1),r}) = 0.} \qquad \forall k \quad \forall m$$

## DeC and residual distribution

Deferred Correction + Residual distribution

- Residual distribution (FV ⇒ FE) ⇒ High order in space
- Prediction/correction/iterations ⇒ High order in time
- Subtimesteps ⇒ High order in time

$$U_\xi^{m,(k+1)} = U_\xi^{m,(k)} - |C_p|^{-1} \sum_{E|\xi \in E} \left( \int_E \Phi_\xi \left( U^{m,(k)} - U^{n,0} \right) d\mathbf{x} + \Delta t \sum_{r=0}^{M} \theta_r^m \mathcal{R}_\xi^E(U^{r,(k)}) \right),$$

with

$$\sum_{\xi \in E} \mathcal{R}_\xi^E(u) = \int_E \nabla_\mathbf{x} F(u) d\mathbf{x}.$$

- The $\mathcal{L}^2$ operator contains also the complications of the spatial discretization (e.g. mass matrix)
- $\mathcal{L}^1$ operator further simplified up to a first order approximation (e.g. **mass lumping**)

# $\mathcal{L}^1$ with mass lumping

Define $\mathcal{L}^1$ as

$$\mathcal{L}^1(\boldsymbol{u}^0, \ldots, \boldsymbol{u}^M) = \begin{pmatrix} \boldsymbol{u}^M - \boldsymbol{u}^0 - \Delta t \beta^M F(\boldsymbol{u}^0) \\ \vdots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \Delta t \beta^1 F(\boldsymbol{u}^0) \end{pmatrix}$$

Define $\mathcal{L}^1$ as

$$\mathcal{L}^1(\boldsymbol{u}^0, \ldots, \boldsymbol{u}^M) = \begin{pmatrix} \boldsymbol{u}^M - \boldsymbol{u}^0 - \Delta t \beta^M \left( F(\boldsymbol{u}^0) + \partial_u F(\boldsymbol{u}^0)(\boldsymbol{u}^M - \boldsymbol{u}^0) \right) \\ \vdots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \Delta t \beta^1 \left( F(\boldsymbol{u}^0) + \partial_u F(\boldsymbol{u}^0)(\boldsymbol{u}^1 - \boldsymbol{u}^0) \right) \end{pmatrix}$$

$$= \begin{pmatrix} \boldsymbol{u}^M - \boldsymbol{u}^0 - \Delta t \beta^M \partial_u F(\boldsymbol{u}^0) \boldsymbol{u}^M \\ \vdots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \Delta t \beta^1 \partial_u F(\boldsymbol{u}^0) \boldsymbol{u}^1 \end{pmatrix}$$

$$\left( I - \Delta t \beta^n \partial_u F(\boldsymbol{u}^0) \right) \boldsymbol{u}^n$$

$$\mathcal{L}^{1,m}(\boldsymbol{u}^0, \dots, \boldsymbol{u}^M) = \boldsymbol{u}^m - \boldsymbol{u}^0 - \Delta t \beta^m \partial_{\boldsymbol{u}} F(\boldsymbol{u}^0) \boldsymbol{u}^m$$

$$\mathcal{L}^{2,m}(\boldsymbol{u}^0, \dots, \boldsymbol{u}^M) = \boldsymbol{u}^m - \boldsymbol{u}^0 - \Delta t \sum_r \theta_r^m F(\boldsymbol{u}^r)$$

ITERATION PROCESS

$$\mathcal{L}^{1,m}(U^{(K)}) - \mathcal{L}^{1,m}(U^{(k-1)}) + \mathcal{L}^{2,m}(U^{(k-1)}) = 0$$

$$U^{m(K)} - U^0 - \Delta t \beta^m \partial_u F(U^0) \cdot U^{m(K)} - U^{m(k-1)} + U^0 + \Delta t \beta^m \partial_u F(U^0) U^{m(k-1)}$$

$$+ U^{m(k)} - U^0 - \Delta t \theta_m^m F(U^r) = 0$$

$$\left[ I - \Delta t \beta^m \partial_u F(U^0) \right] \left( U^{m(K)} - U^{m(k-1)} \right) = \mathcal{L}^{2,m}(U^{(k-1)}) \qquad K = 1 \to \bar{K}$$

$$\forall m$$

D. Torlo    ADER vs DeC

$$\boldsymbol{u}^{(k),m} - \boldsymbol{u}^0 - \Delta t \sum_{r=0}^{M} \theta_r^m F(\boldsymbol{u}^{(k-1),r}) = 0$$

$$\forall k = 1, \underline{\quad} \overline{K}$$

$$\forall m = 1, \underline{\quad}, M$$

$$\underline{U}^{(1)} = \begin{pmatrix} U^{(1),1} \\ U^{(1),n} \end{pmatrix}$$

$$\underline{U}^{(2)}$$

$$ST \\ = K \cdot n$$

We can write DeC as RK defining $\underline{\theta}_0 = \{\theta_0^m\}_{m=1}^M$, $\underline{\theta}^M = \theta_r^M$ with $r \in 1, \ldots, M$, denoting the vector $\underline{\theta}_r^{M,T} = (\theta_1^M, \ldots, \theta_M^M)$. The Butcher tableau for an arbitrarily high order DeC approach is given by:

$$
\begin{array}{c|ccccccc}
0 & 0 \\
\beta & \beta \\
\beta & \underline{\theta}_0 & \underline{\tilde{\theta}} \\
\vdots & \underline{\theta}_0 & \underline{0} & \underline{\tilde{\theta}} \\
\vdots & \underline{\theta}_0 & \underline{0} & \underline{0} & \underline{\tilde{\theta}} \\
\vdots & \vdots & \vdots & \vdots & \ddots & \ddots \\
\beta & \underline{\theta}_0 & \underline{0} & \cdots & \cdots & \underline{0} & \underline{\tilde{\theta}} \\
\hline
 & \theta_0^M & \underline{0}^T & \cdots & & \cdots & \underline{0}^T & \underline{\theta}_r^{M,T}
\end{array}
\tag{6}
$$

Idea: study the RK version!

$$\bullet \boxed{u' = \lambda u} \quad \boxed{\Re(\lambda) < 0.} \tag{7}$$

$$\underbrace{u_{n+1} = \underbrace{R(\lambda \Delta t)}u_n,} \quad \underbrace{R(z) = \underbrace{1 + z\underline{b}^T (I - zA)^{-1}\mathbf{1}}}, \quad z = \lambda \Delta t \tag{8}$$

Goal: find $z \in \mathbb{C}$ such that $|R(z)| < 1$.

Recall: stability function for explicit RK methods is a polynomial, indeed the inverse of $(I - zA)$ can be written in Taylor expansion as

$$\underbrace{(I - zA)^{-1} = \sum_{r=0}^{\infty} z^r A^s = I + zA + z^2 A^2 + \ldots,} \tag{9}$$

and, since $A$ is strictly lower triangular, it is nilpotent. Hence, $R(z)$ is a polynomial in $z$ with degree at most equal to $S$.



$$A^s = \mathcal{O}$$

## Stability of (explicit) DeC

**Theorem** HAILER BECK

*If the RK method is of order P, then*

$$R(z) = \underbrace{1 + z + \frac{z^2}{2!} + \cdots + \frac{z^P}{P!}} + \underbrace{O(z^{P+1})}. \tag{10}$$

The first $P + 1$ terms of the stability functions $R(\cdot)$ for explicit DeCs of order $P$ are known.

**Theorem**

*The stability function of any explicit DeC of order $P$ (with $P$ iterations) is*

$$R(z) = \sum_{r=0}^{P} \frac{z^r}{r!} = \underbrace{1 + z + \frac{z^2}{2!} + \cdots + \frac{z^P}{P!}} \tag{11}$$

*and does not depend on the distribution of the subtimenodes.*

## Proof (1/3)

$$A = \begin{pmatrix} 0 & 0 & 0 & \ldots & 0 & 0 \\ \star & 0 & 0 & \ldots & 0 & 0 \\ \star & \star & 0 & \ldots & 0 & 0 \\ \star & 0 & \star & \ldots & 0 & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ \star & 0 & 0 & \ldots & \star & 0 \end{pmatrix},$$

Block structure of the matrix $A$

$\star$ are some non-zero block matrices and the 0 are some zero block matrices.

The number of blocks in each line and row of these matrices is $P$, the order of the scheme.

$$A^P = 0$$

## Proof (2/3)

By induction, $A^k$ has zeros in the upper triangular part, in the main block diagonal, and in all the $k-1$ block diagonals below the main diagonal, i.e.,

$$(A^k)_{i,j} = 0 \quad , \text{if } i < j + k,$$

where the indexes here refer to the blocks. Indeed, it is true that $A_{i,j} = 0$ if $i < j + 1$. Now, let us consider the entry $(A^{k+1})_{i,j}$ with $i < j + k + 1$, i.e., $i - k < j + 1$. It is defined as

$$(A^{k+1})_{i,j} = \sum_w (A^k)_{i,w} A_{w,j}. \tag{12}$$

Now, we can prove that all the terms of the sum are 0. Let $w < j + 1$, then $A_{w,j} = 0$ because of the structure of $A$; while, if $w \geq j + 1 > i - k$, we have that $i < w + k$, so $(A^k)_{i,w} = 0$ by induction.

## Proof (3/3)

In particular, this means that $A^P = \underline{0}$, because $i$ is always smaller than $j + P$ as $P$ is the number of the block matrices that we have. Hence,

$$(I - zA)^{-1} = \sum_{r=0}^{\infty} z^r A^s = \sum_{r=0}^{P-1} z^r A^s = I + zA + z^2 A^2 + \cdots + z^{P-1} A^{P-1}. \tag{13}$$

Plugging this result into $R(z) = 1 + zb^T (I - zA)^{-1} \mathbf{1}$, the stability function $R(z)$ is a polynomial of degree $P$, the order of the scheme. All terms of order lower or equal to $P$ must agree with the expansion of the exponential function, so it must be

$$R(z) = \sum_{r=0}^{P} \frac{z^r}{r!} = 1 + z + \frac{z^2}{2!} + \cdots + \frac{z^P}{P!}. \tag{14}$$

Note: no assumption on the distribution of the subtimenodes.

## CODE

- Choice of iterations ($P$) and order
- Choice of point distributions $t^0, \ldots, t^M$
- Computation of $\theta$
- Loop for timesteps
- Loop for correction
- Loop for subtimesteps

## Outline

- Cauchy–Kovalevskaya theorem
- Modern automatic version *2008*
  - Space/time DG
  - Prediction/Correction
  - Fixed-point iteration process

Prediction: iterative procedure

Modern approach is DG in space time for hyperbolic problem    *FV*    *FEM*

$$\partial_t u(x,t) + \nabla \cdot F(u(x,t)) = 0, \quad x \in \Omega \subset \mathbb{R}^d,\ t > 0. \quad (15)$$

*WEAK FORMULATION*

$$\sum_{p,q} \int_{T^n \times V_i} \overbrace{\theta_{rs}(x,t)}\underbrace{\partial_t \theta_{pq}(x,t)z^{pq}}\mathrm{d}x\mathrm{d}t + \int_{T^n \times V_i} \overbrace{\theta_{rs}(x,t)}\nabla_{\mathbf{x}} \cdot F(\underbrace{\theta_{pq}(x,t)z^{pq}}_{RECONSTR})\mathrm{d}x\mathrm{d}t = 0. \quad \forall r,s$$

$\theta_{rs}(x,t)$   *TEST*

Correction step: communication between cells

$$\int_{V_i} \Phi_r \left( u(t^{n+1}) - u(t^n) \right)\mathrm{d}x + \int_{T^n \times \partial V_i} \Phi_r(x)\overbrace{\mathcal{G}(\underline{z^-}, \underline{z^+})} \cdot \mathbf{n}\,\mathrm{d}S\,\mathrm{d}t - \underbrace{\int_{T^n \times V_i} \nabla_{\mathbf{x}}\Phi_r \cdot F(z)\,\mathrm{d}x\,\mathrm{d}t}_{} = 0,$$

Defining $\theta_{rs}(x, t) = \Phi_r(x)\phi_s(t)$ basis functions in space and time

$$\int_{T^n \times V_i} \theta_{rs}(x, t)\partial_t \theta_{pq}(x, t)u^{pq}\mathrm{d}x\mathrm{d}t + \int_{T^n \times V_i} \theta_{rs}(x, t)\nabla \cdot F(\theta_{pq}(x, t)u^{pq})\mathrm{d}x\mathrm{d}t = 0. \tag{16}$$

## ADER: space-time discretization

Defining $\theta_{rs}(x, t) = \Phi_r(x)\phi_s(t)$ basis functions in space and time

$$\underbrace{\int_{T^n \times V_i} \theta_{rs}(x, t) \partial_t \theta_{pq}(x, t) u^{pq} \mathrm{d}x\mathrm{d}t}_{} + \underbrace{\int_{T^n \times V_i} \theta_{rs}(x, t) \nabla \cdot F(\theta_{pq}(x, t) u^{pq}) \mathrm{d}x\mathrm{d}t = 0.} \tag{16}$$

This leads to

$$\underline{\underline{\mathrm{M}}}_{rspq} u^{pq} = \underline{r}(\underline{\underline{u}})_{rs}, \tag{17}$$

solved with fixed point iteration method.

+ Correction step where cells communication is allowed (derived from (16)).

$\text{① } DE \qquad \partial_t u = F(u)$

Simplify! Take $\boldsymbol{u}(t) = \sum_{m=0}^{M} \underbrace{\phi_m(t)}\boldsymbol{u}^m = \underline{\phi}(t)^T \underline{\boldsymbol{u}}$

WEAK FORM OF ODE $\quad \int_{T^n} \widetilde{\psi(t)} \partial_t \boldsymbol{u}(t)\,dt - \int_{T^n} \psi(t) F(\boldsymbol{u}(t))\,dt = 0, \quad \forall \psi : T^n = \underbrace{[t^n, t^{n+1}]} \to \mathbb{R}.$

$\underline{\mathcal{L}^2(\underline{\boldsymbol{u}})} := \int_{T^n} \underline{\phi}(t) \partial_t \underline{\phi}(t)^T \underline{\boldsymbol{u}}\,dt - \int_{T^n} \underline{\phi}(t) F(\underline{\phi}(t)^T \underline{\boldsymbol{u}})\,dt = 0$

$\underline{\phi}(t) = (\phi_0(t), \ldots, \phi_M(t))^T$

$\boxed{\int_{T^n} \phi_i^{(t)} \partial_t \phi_j^{(t)} u^J \, dt} - \int_{T^n} \phi_i \, F(\phi_j(t) u^J) dt$

Quadrature...

INTEGRATION BY PARTS

$= \phi_i(t^{n+1}) \phi_j(t^{n+1}) u^J - \phi_i(t^n) \phi_j(t^n) u^J$

$\underbrace{\qquad}_{U(t^n)}$

$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\,\underline{\boldsymbol{u}} - \underline{r}(\underline{\boldsymbol{u}}) = 0 \iff \underline{\underline{M}}\,\underline{\boldsymbol{u}} = \underline{r}(\underline{\boldsymbol{u}}).$ (18)

$- \int_{t^n}^{t^{n+1}} \partial_t \phi_i(t) \cdot \phi_j(t) u^J \, dt$

Nonlinear system of $M \times S$ equations

$\phi_i(t^n) \phi_j(t^{n+1}) u^J - \int_{t^n}^{t^{n+1}} \phi_i'(t) \phi_j(t) \, dt \cdot u^J$

$- \int_{T^n} \phi_i \, F(\phi_j \, u^J) \, dt$

$- \underbrace{\phi_i(t^n) U(t^n)}_{\text{EXPLICIT}} - \int_{t^n}^{t^{n+1}} \phi_i \, F(\phi_j(t) u^J) \, dt = \mathcal{L}^{2,i}$

# ADER: Mass matrix

What goes into the mass matrix? Use of the integration by parts

WEAK FERMULATION

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \int_{T^n} \underline{\phi}(t)\partial_t\underline{\phi}(t)^T\underline{\boldsymbol{u}}\,dt + \int_{T^n} \underline{\phi}(t)F(\underline{\phi}(t)^T\underline{\boldsymbol{u}})\,dt =$$

I.B.P. IN TIME $\Rightarrow$
$$\underline{\phi}(t^{n+1})\underline{\phi}(t^{n+1})^T\underline{\boldsymbol{u}} - \underline{\phi}(t^n)\boldsymbol{u}^n - \int_{T^n} \partial_t\underline{\phi}(t)\underline{\phi}(t)^T\underline{\boldsymbol{u}} - \int_{T^n} \underline{\phi}(t)F(\underline{\phi}(t)^T\underline{\boldsymbol{u}})\,dt$$

KNOWN

$$\underline{\underline{M}} = \underline{\phi}(t^{n+1})\underline{\phi}(t^{n+1})^T - \int_{T^n} \partial_t\underline{\phi}(t)\underline{\phi}(t)^T$$

$$\underline{r}(\underline{\boldsymbol{u}}) = \underline{\phi}(t^n)\boldsymbol{u}^n + \int_{T^n} \underline{\phi}(t)F(\underline{\phi}(t)^T\underline{\boldsymbol{u}})\,dt$$

$$\underline{\underline{M}}\underline{\boldsymbol{u}} = \underline{r}(\underline{\boldsymbol{u}})$$

NON LINEAR SYS

LINEAR    NONLINEAR    DIM $(M \times S)$

Iterative procedure to solve the problem for each time step

$$\underline{u}^{(k)} = \underline{\underline{M}}^{-1}\underline{r}(\underline{u}^{(k-1)}), \quad k = 1, \ldots, \text{convergence} \tag{19}$$

with $\underline{u}^{(0)} = u(t^n)$.
Reconstruction step

$$u(t^{n+1}) = u(t^n) - \int_{T^n} F(u^{(K)}(t))dt.$$

- Convergence?
- How many steps $K$?
- Accuracy $\mathcal{L}^2$?  $\longrightarrow$ IMPLICIT RK

## ADER 2nd order

Example with 2 Gauss Legendre points, Lagrange polynomials and 2 iterations
Let us consider the timestep interval $[t^n, t^{n+1}]$, rescaled to $[0,1]$.
Gauss-Legendre points quadrature and interpolation (in the interval $[0,1]$)

$$\underline{t}_q = \left( t_q^0, t_q^1 \right) = \left( t^0, t^1 \right) = \left( \frac{\sqrt{3}-1}{2\sqrt{3}}, \frac{\sqrt{3}+1}{2\sqrt{3}} \right), \quad \underline{w} = (1/2, 1/2).$$

$$\underline{\phi}(t) = (\phi_0(t), \phi_1(t)) = \left( \frac{t - t^1}{t^0 - t^1}, \frac{t - t^0}{t^1 - t^0} \right).$$

Then, the mass matrix is given by

$$\underline{\underline{M}}_{m,l} = \phi_m(1)\phi_l(1) - \phi'_m(t^l)w_l, \quad m, l = 0, 1,$$

$$\underline{\underline{M}} = \begin{pmatrix} 1 & \frac{\sqrt{3}-1}{2} \\ -\frac{\sqrt{3}+1}{2} & 1 \end{pmatrix}.$$

*(handwritten annotations)*

Top right:
$$\begin{array}{c|c} \text{POLY} & \text{QUAD} \\ \text{LAG} & \\ \hline \text{GAUSS} & \text{GAUSS} \\ \text{LOBATTO} & \text{LOBATTO} \\ \text{X EGV1} & \text{GAUSS} \end{array}$$

$p$ points

$\phi_n(t^{n+1}) \phi_l(t^{n+1}) - \int_{t^n}^{t^{n+1}} \phi'_n(t) \cdot \phi_l(t) \, dt$

$\mathbb{P}^{P-2}$

$-\int_0^1 \phi'_m \, \phi_J \, dt$

$QUAD = LAGR \cdot POLY$

$- w_l \cdot \phi'_m(t^l)$

$$\frac{2}{3}\phi_J \, U^3$$

The right hand side is given

$$\phi_m \, u^m + \Delta t \int_0^1 \phi_m \, F(u(t)) \, dt \simeq$$

$$r(\underline{u})_m = u(0)\phi_m(0) + \underline{\Delta t} F(u(t^m)) w_m, \quad m = 0, 1.$$

$$\overline{\phi_m \cdot U_J \Delta t \, w_m \cdot F(u^m)}$$

$$\phi_J(t^m) = \delta_{Jm}$$

$$\underline{r}(\underline{u}) = u(0)\underline{\phi}(0) + \Delta t \begin{pmatrix} F(u(t^1)) w_1 \\ F(u(t^2)) w_2. \end{pmatrix}.$$

Then, the coefficients $\underline{u}$ are given by

$$\boxed{\underline{u}^{(k+1)} = \underline{\underline{M}}^{-1} \underline{r}(\underline{u}^{(k)}).}$$

$$= \overbrace{(M^{-1} \phi(1)^T \, U^m}^{= U^m}$$

$$+ \Delta t \, \underline{\underline{M}}^{-1} \underline{\underline{R}} \, F(\underline{u}^{(k)})}_{\text{eval } m\pi}$$

$$R_{HS}$$

Finally, use $\underline{u}^{(k+1)}$ to reconstruct the solution at the time step $t^{n+1}$:

$$\sum \phi_J(1) \, U^{J,(k+1)}$$

$$\underline{u}^{n+1} = \underline{\phi}(1)^T \underline{u}^{(k+1)} = u^n + \int_{T^n} \underline{\phi}(t)^T dt \, F(\underline{u}^{(k)}).$$

## CODE

- Choice: $\phi$ Lagrangian basis functions $(TRIAL = TEST)$
- Different subtimesteps: Gauss-Legendre, Gauss–Lobatto, equispaced
  - $auto | GAUSS$   $LOB$   $GAUSS$
- Precompute $\underline{\underline{M}}$
- Precompute the rhs vector part using quadratures after a further approximation

$$\underline{r}(\underline{u}) = \underbrace{\phi(t^n)}\underline{u}^n + \int_{T^n} \underline{\phi}(t) F(\underline{\phi}(t)^T \underline{u}) dt \approx \underline{\phi}(t^n)\underline{u}^n + \overbrace{\int_{T^n} \underline{\phi}(t)\underline{\phi}(t)^T dt}\, F(\underline{u})$$

$$\underline{\phi}(0)$$

$\underbrace{\qquad}_{\text{Can be stored}}$

- Precompute the reconstruction coefficients $\underline{\phi}(1)^T$

$$\underline{\underline{R}}_{ij} = \int \varphi_i \varphi_j \, dt$$

$$IF \ QUAD \qquad \simeq W_i \, \delta_{ij}$$
$$= LAGRA \ POINT$$

# Outline

## ADER[6] and DeC[7]: immediate similarities

ARBITRARILY

- High order time(space) discretization

- Start from a well known space discretization (FE/DG/FV)

- FE reconstruction in time

- System in time, with $M$ equations        $\mathcal{L}^2 \geq 0$

- Iterative method / $K$ corrections

---

[6]M. Dumbser, D. S. Balsara, E. F. Toro, and C.-D. Munz. A unified framework for the construction of one-step finite volume and discontinuous galerkin schemes on unstructured meshes. Journal of Computational Physics, 227(18):8209–8253, 2008.

[7]R. Abgrall. High order schemes for hyperbolic problems using globally continuous approximation and avoiding mass matrices. Journal of Scientific Computing, 73(2):461–494, Dec 2017.

## ADER[6] and DeC[7]: immediate similarities

- High order time-space discretization
- Start from a well known space discretization (FE/DG/FV)
- FE reconstruction in time
- System in time, with $M$ equations
- Iterative method / $K$ corrections

- Both high order explicit time integration methods (neglecting spatial discretization)

---

[6]M. Dumbser, D. S. Balsara, E. F. Toro, and C.-D. Munz. A unified framework for the construction of one-step finite volume and discontinuous galerkin schemes on unstructured meshes. Journal of Computational Physics, 227(18):8209–8253, 2008.

[7]R. Abgrall. High order schemes for hyperbolic problems using globally continuous approximation and avoiding mass matrices. Journal of Scientific Computing, 73(2):461–494, Dec 2017.

ADER $\qquad \Pi \, u^{(K)} = \mathcal{R}(u^{(K-1)})$ $\overset{\text{GOAL}}{\Longleftarrow}$

DeC $\qquad \mathcal{L}^{1}(\hat{u}) = U^{m} - U^{0} + \Delta t \bar{F}(U^{0}) \cdot \beta^{m}$

$\mathcal{L}^{2,m}(U) = U^{m} - U^{0} + \Delta t \sum_{r=0}^{M} \Theta_{r}^{m} F(U^{r})$

$\begin{cases} \mathcal{L}^{2}(\underline{u}) = \Pi \, \underline{u} - \mathcal{R}(\underline{u}) \\ \mathcal{L}^{2}(\underline{u}) = \Pi \, \underline{u} - \underbrace{\mathcal{R}(\underline{u}^{m})}_{\text{explicit}} \end{cases}$ $\begin{array}{l} \text{FIRST} \\ \text{APPROXIMATION} \end{array}$

$\mathcal{L}^{1}(U^{(K+1)}) = \mathcal{L}^{1}(U^{(K)}) - \mathcal{L}^{2}(U^{(K)})$

$\mathcal{L}^{1}(\underline{u}^{(K)}) = \mathcal{L}^{1}(u^{(K-1)}) - \mathcal{L}^{2}(u^{(K-1)})$

$\underline{\underline{\Pi}} \, \underline{u}^{(K)} - \cancel{\mathcal{R}(\underline{u}^{m})} = \cancel{\Pi \, \underline{u}^{(K-1)}} - \mathcal{R}(\underline{u}^{m}) - \left( \cancel{\Pi \, \underline{u}^{(K-1)}} - \mathcal{R}(\underline{u}^{(K-1)}) \right)$

$\underline{\underline{\Pi}} \, \underline{u}^{(K)} = \mathcal{R}(\underline{u}^{(K-1)}) \qquad \checkmark$

$\checkmark \mathcal{L}^2$ IS COERCIVE $\qquad \checkmark \mathcal{L}^2 - \mathcal{L}^1$ IS LIPSCHITZ CONT. Cnst $\Delta t \cdot C_2$

· $\exists^1_{\underset{u^*}{x}} \mathcal{L}^2(v^*) = 0 \qquad \Rightarrow \qquad || u^{(K)} - v^* || \leq (c \Delta t)^K || v^{(0)} - v^* ||$

$(C) || \mathcal{L}^2(u) - \mathcal{L}^2(v) || \geq C_1 || u - v ||$

$$\begin{cases} \mathcal{L}^2(\underline{u}) = \Pi \, \underline{u} \, - \mathcal{R}(\underline{u}) \\ \mathcal{L}^2(\underline{u}) = \Pi \, \underline{u} \, - \mathcal{R}(u^{\wedge}) \end{cases}$$

$|| \underline{\underline{\Pi}} \, \underline{u} - \mathcal{R}(\underline{u^{\wedge}}) - \underline{\underline{\Pi}} \, \underline{v} - \mathcal{R}(\underline{u^{\wedge}}) ||$

$= || \underline{\underline{\Pi}} \, (\underline{u} - \underline{v}) || > C_1 (\underline{\underline{\Pi}}) \, || \underline{u} - \underline{v} ||$

$\qquad \qquad (|| \pi^{-1} ||)$

$(L.) \, || \mathcal{L}^1(u) - \mathcal{L}^2(u) - \mathcal{L}^1(v) + \mathcal{L}^2(v) || = || \mathcal{R}(\underline{u}) - \mathcal{R}(v^{\wedge}) - \mathcal{R}(\underline{v}) + \mathcal{R}(v^{\wedge}) ||$

$\qquad = || \mathcal{R}(u) - \mathcal{R}(v) || \leq || u^{\vee} + \Delta t \underline{\underline{R}} F(\underline{u}) - u^{\wedge} + \Delta t \underline{\underline{R}} F(\underline{v}) || \leq \underline{\Delta t} \cdot \underbrace{|| \underline{\underline{R}} || \, L}_{C_2} \, || \underline{u} - \underline{v} ||$

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\underline{\boldsymbol{u}}),$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\boldsymbol{u}(t^n)).$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

$$\underline{\underline{M}}\boldsymbol{u}^{(k)} - r(\boldsymbol{u}^{(k),0}) - \underline{\underline{M}}\boldsymbol{u}^{(k-1)} + r(\boldsymbol{u}^{(k-1),0}) + \underline{\underline{M}}\boldsymbol{u}^{(k-1)} - r(\underline{\boldsymbol{u}}^{(k-1)}) = 0$$

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\underline{\boldsymbol{u}}),$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\boldsymbol{u}(t^n)).$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

$$\underline{\underline{M}}\boldsymbol{u}^{(k)} - \cancel{r(\boldsymbol{u}^{(k),0})} - \underline{\underline{M}}\boldsymbol{u}^{(k-1)} + \cancel{r(\boldsymbol{u}^{(k-1),0})} + \underline{\underline{M}}\boldsymbol{u}^{(k-1)} - r(\underline{\boldsymbol{u}}^{(k-1)}) = 0$$

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{\mathrm{M}}}\underline{\boldsymbol{u}} - r(\underline{\boldsymbol{u}}),$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}) := \underline{\underline{\mathrm{M}}}\underline{\boldsymbol{u}} - r(\boldsymbol{u}(t^n)).$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

$$\underline{\underline{\mathrm{M}}}\boldsymbol{u}^{(k)} - \cancel{r(\boldsymbol{u}^{(k),0})} - \cancel{\underline{\underline{\mathrm{M}}}\boldsymbol{u}^{(k-1)}} + \cancel{r(\boldsymbol{u}^{(k-1),0})} + \cancel{\underline{\underline{\mathrm{M}}}\boldsymbol{u}^{(k-1)}} - r(\underline{\boldsymbol{u}}^{(k-1)}) = 0$$

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\underline{\boldsymbol{u}}),$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\boldsymbol{u}(t^n)).$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(k)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(k-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(k-1)}), \qquad k = 1, \ldots, K,$$

$$\underline{\underline{M}}\boldsymbol{u}^{(k)} - \cancel{r(\boldsymbol{u}^{(k),0})} - \cancel{\underline{\underline{M}}\boldsymbol{u}^{(k-1)}} + \cancel{r(\boldsymbol{u}^{(k-1),0})} + \cancel{\underline{\underline{M}}\boldsymbol{u}^{(k-1)}} - r(\underline{\boldsymbol{u}}^{(k-1)}) = 0$$
$$\underline{\underline{M}}\boldsymbol{u}^{(k)} - r(\underline{\boldsymbol{u}}^{(k-1)}) = 0.$$

$$\mathcal{L}^2(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\underline{\boldsymbol{u}}),$$
$$\mathcal{L}^1(\underline{\boldsymbol{u}}) := \underline{\underline{M}}\underline{\boldsymbol{u}} - r(\boldsymbol{u}(t^n)).$$

Apply the DeC Convergence theorem!

- $\mathcal{L}^1$ is coercive because $\underline{\underline{M}}$ is always invertible
- $\mathcal{L}^1 - \mathcal{L}^2$ is Lipschitz with constant $C\Delta t$ because they are consistent approx of the same problem
- Hence, after $K$ iterations we obtain a $K$th order accurate approximation of $\underline{\boldsymbol{u}}^*$

$$\| u^{(K)} - v^* \| \leq (C \Delta t)^K \| u^{(0)} - v^* \|$$

$$\mathcal{L}^2(\boldsymbol{u}^0, \dots, \boldsymbol{u}^M) := \begin{cases} \boldsymbol{u}^M - \boldsymbol{u}^0 - \sum_{r=0}^{M} \int_{t^0}^{t^M} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \\ \dots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \sum_{r=0}^{M} \int_{t^0}^{t^1} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \end{cases}.$$

? INTO
WEAK FORM ?

$$\mathcal{L}^2(\boldsymbol{u}^0, \dots, \boldsymbol{u}^M) := \begin{cases} \boldsymbol{u}^M - \boldsymbol{u}^0 - \sum_{r=0}^M \int_{t^0}^{t^M} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \\ \dots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \sum_{r=0}^M \int_{t^0}^{t^1} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \end{cases} .$$

$$\mathcal{L}^2(\boldsymbol{u}^0, \ldots, \boldsymbol{u}^M) := \begin{cases} \boldsymbol{u}^M - \boldsymbol{u}^0 - \sum_{r=0}^M \int_{t^0}^{t^M} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \\ \ldots \\ \boldsymbol{u}^1 - \boldsymbol{u}^0 - \sum_{r=0}^M \int_{t^0}^{t^1} F(\boldsymbol{u}^r)\varphi_r(s)\mathrm{d}s \end{cases}.$$

$$\underbrace{\chi_{[t^0,t^m]}(t^m)\boldsymbol{u}^m}_{} - \underbrace{\chi_{[t^0,t^m]}(t_0)\boldsymbol{u}^0}_{} - \int_{t^0}^{\overbrace{t^m}} \underbrace{\chi_{[t^0,t^m]}(t)}_{} \sum_{r=0}^M F(\boldsymbol{u}^r)\varphi_r(t)\mathrm{d}t = 0$$

$$\int_{t^0}^{t^M} \chi_{[t^0,t^m]}(t)\partial_t\left(\boldsymbol{u}(t)\right)\mathrm{d}t - \int_{t^0}^{\overbrace{t^M}} \underbrace{\chi_{[t^0,t^m]}(t)}_{} \sum_{r=0}^M F(\boldsymbol{u}^r)\varphi_r(t)\mathrm{d}t = 0,$$

$$\int_{T^n} \underbrace{\psi_m(t)}_{}\partial_t\boldsymbol{u}(t)\mathrm{d}t - \underbrace{\int_{T^n} \psi_m(t)F(\boldsymbol{u}(t))\mathrm{d}t = 0.}_{}$$

## Classical Runge Kutta (RK)

- One step method
- Internal stages

Explicit Runge Kutta

+ Simple to code
- Not easily generalizable to arbitrary order
- Stages $>$ order

Implicit Runge Kutta

+ Arbitrarily high order
- Require nonlinear solvers for nonlinear systems
- May not converge

## DeC – ADER

- One step method
- Internal subtimesteps  + ITERATIONS
- Can be rewritten as explicit RK (for ODE)
+ Explicit
+ Simple to code
+ Iterations = order
+ Arbitrarily high order
- Large memory storage

# Outline

## Stability

Since ADER can be written as a DeC, the stability functions are given by the same formula as for DeC and the stability regions are the following.



Figure: Stability region

## Accuracy of ADER $\mathcal{L}^2$ operators $\quad \mathcal{L}^2 = \underline{\Pi} \, \underline{u} - \underline{\Pi}(\underline{u}) \quad ?$ ORDER ( IMPLICIT RK)

The two things that determine the accuracy of the ADER method are the iterations $P$ and the accuracy of $\mathcal{L}^2$.

### Accuracy of ADER $\mathcal{L}^2$ for different distributions

- Equispaced: boring, minimum accuracy possible $M + 1$ nodes $p = M + 1$
- Guass–Lobatto: this generates the LobattoIIIC methods, $M + 1$ nodes $p = 2M$, S styes 2S-2 order
- Gauss–Legendre: this does not generate Gauss methods, $M + 1$ nodes $p = 2M + 1$, S stages 2S-1 order

ADER EXPLICIT

GLG     $M + 1$ NODES   $\Rightarrow \mathcal{L}^2$ ORDER $2M + 1$   $\Rightarrow K = 2M + 1$  ✓

GLB     $M + 1$ NODES   $\Rightarrow \mathcal{Y}^2$   $2M$   $\Rightarrow K = 2M$  ✓

$\Rightarrow P = K$

EINSTEIN NOTATION

Here, we see $\mathcal{L}^2$ as an implicit RK

$$\mathcal{L}^{2,m}(\underline{\boldsymbol{u}}) = \underset{j}{\underline{\underline{\mathrm{M}}}}^m_j \boldsymbol{u}^{(j)} - \underline{\phi}^m(t^n)\boldsymbol{u}^n - \underbrace{\int_{T^n} \underline{\phi}^m(t)\underline{\phi}(t)_j dt}_{\Delta t \underline{\underline{\mathrm{R}}}^m_j} F(\boldsymbol{u}^{(j)}) = 0$$

$$\underbrace{(\Pi^{-1})^z}_{} \qquad \underbrace{\qquad}_{1}$$

$$\tilde{\mathcal{L}}^{2,z}(\underline{\boldsymbol{u}}) = \boldsymbol{u}^{(z)} - \underbrace{(\underline{\underline{\mathrm{M}}}^{-1})^z_m \underline{\phi}^m(t^n)}_{}\boldsymbol{u}^n - \Delta t \underbrace{(\underline{\underline{\mathrm{M}}}^{-1})^z_m \underline{\underline{\mathrm{R}}}^m_j}_{} F(\boldsymbol{u}^{(j)}) = 0$$

$RK$ $\qquad \boldsymbol{u}^{(z)} = \boldsymbol{u}^n + \Delta t\, a_{z,j} F(\boldsymbol{u}^{(j)})$

- $a_{mj} := (\underline{\underline{\mathrm{M}}}^{-1})^z_m \underline{\underline{\mathrm{R}}}^m_j$ ✓

$$u^{\wedge n} = u^{\wedge} + \frac{1}{\Delta t} \underbrace{\int_0^1 \phi_r(t)\, dt}_{= w_r = b_r} \cdot F(u^r)$$

- Prove that $(\underline{\underline{\mathrm{M}}}^{-1})^z_m \underline{\phi}^m(t^n) = 1$ for every $z$ ✓
- $c^m = \sum_r a_{mr} = t^m$ ?
- $b_r = \frac{1}{\Delta t}\int_{T^m} \phi_r(t) dt = w_r$ quadrature weights ✓

- $(\Pi^{-1})^z_{\,J} \, \phi^J(0) \stackrel{?}{=} \mathbb{1}^z \quad \forall z = \overset{t^o}{0}, \; - \; , \overset{t^n}{n}$

$$\Longleftrightarrow \quad \Pi^m_{\;z} \, (\Pi^{-1})^z_{\,J} \, \phi^J(0) \stackrel{?}{=} \Pi^m_{\;z} \, \mathbb{1}^z \quad \Longleftrightarrow \quad \phi^m(0) \stackrel{?}{=} \Pi^m_{\;z} \cdot \mathbb{1}^z$$

$$\underbrace{\phantom{\Pi^m_{\;z} (\Pi^{-1})^z_{\,J}}}_{\mathcal{I}^m_{\,J}}$$

$$\sum_{z} \Pi^m_{\;z} \cdot \mathbb{1}^z = \sum_{z=0}^{n} \Pi^m_{\;z} = \sum_{z=0}^{n} \phi_m(1) \, \phi_z(1) - \int_0^1 \phi'_m(t) \cdot \phi_z(t) \, dt =$$

$$\overset{\mathbb{P}^{n-1}}{\uparrow} \times \overset{\mathbb{P}^{n}}{\uparrow} \in \mathbb{P}^{2n-1}$$

RECALL THAT GLC, CLB QUADRATURE RULE IS EXACT POLY DEGREE $\underline{2S-3}$

$\qquad\qquad\quad\; 2S-1 \quad 2S-3 \qquad\qquad\qquad\qquad\qquad\qquad\qquad 2(n+1)-3$

RECALL LAGR. BASIS FUNCTION $\sum_{J=0}^{n} \phi_J(t) \equiv \mathbb{1} \quad \forall t \qquad 2n-1$

$$= \phi_m(1) \cdot \mathbb{1} - \int_0^1 \phi'_m(t) \, dt = \phi_m(1) - \big[\phi_m(t)\big]_0^1 = \cancel{\phi_m(1)} - \cancel{\phi_m(1)} + \phi_m(0)$$

$$\overset{\mathbb{P}^{n-1}}{\uparrow} \qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad \boxtimes$$

$\sum_k a_{zk} = c_z = t^z$

2) $\sum_k (\Pi^{-1})^{\mathrm{J}}_{\mathrm{J}} R^{\mathrm{J}}_K \mathbb{1}^K \overset{?}{=} t^z \quad \hookleftarrow$

$R^{\mathrm{J}}_k \cdot \mathbb{1}^K \overset{?}{=} \Pi^{\mathrm{J}}_z \cdot t^z$

$R^{\mathrm{J}}_k \cdot \mathbb{1}^K = \sum_K \int_0^1 \phi_{\mathrm{J}}(t) \underbrace{\phi_K(t)}_{=1} = \int_0^1 \phi_{\mathrm{J}}(t)\, dt = \underline{w_{\mathrm{J}}}$

$\Pi^{\mathrm{J}}_z\, t^z = \phi_{\mathrm{J}}(1)\, \phi_z(1) \cdot t^z - \int_0^1 \phi'_{\mathrm{J}}(t) \cdot \overbrace{\phi_z(t) \cdot t^z}^{} \underset{\substack{t=1 \\ \Downarrow}}{=} \phi_{\mathrm{J}}(1) \cdot \mathbb{1} - \int_0^1 \underbrace{\phi'_{\mathrm{J}}(t) \cdot t\, dt}_{\in \phi^n}$

$\underset{\substack{\uparrow \\ \text{INTERPOLATING } t}}{\sum_z \phi_z(t) \cdot t^z} \underset{\Downarrow}{=} \underline{t} \qquad \text{NODES } t^z \qquad t \in \mathbb{P}$

EXACT INTER

$\overset{\text{I.BP.}}{=} \phi_{\mathrm{J}}(1) - \left[\phi_{\mathrm{J}}(t) \cdot t\right]_0^1 + \int_0^1 \phi_{\mathrm{J}}(t) \cdot \mathbb{1}\, dt = \phi_{\mathrm{J}}(1) - \phi_{\mathrm{J}}(1) - 0 + \underline{w_{\mathrm{J}}}$

EXACT QUAD

$(t)^1$

## BCD conditions (Butcher 1964)

Define the conditions

$B(2S-2)$ GL
$B(2S)$ GL

$$B(p): \quad \sum_{i=1}^{s} b_i c_i^{z-1} = \frac{1}{z}, \qquad\qquad z = 1, \ldots, p; \qquad (20)$$

$\cancel{C(S)}$  $C(S-1)$

$$C(\eta): \quad \sum_{j=1}^{s} a_{ij} c_j^{z-1} = \frac{c_i^z}{z}, \qquad i = 1, \ldots, s, \ z = 1, \ldots, \eta; \qquad (21)$$

$D(S-1)$

$$D(\zeta): \quad \sum_{i=1}^{s} b_i c_i^{z-1} a_{ij} = \frac{b_j}{z}(1 - c_j^z), \qquad j = 1, \ldots, s, \ z = 1, \ldots, \zeta. \qquad (22)$$

### Theorem (Butcher 1964)

*If the coefficients $b_i, c_i, a_{ij}$ of a RK scheme satisfy $B(p)$, $C(\eta)$ and $D(\zeta)$ with $p \leq \eta + \zeta + 1$ and $p \leq 2\eta + 2$, then the method is of order $p$.*

$\underline{C(s-1)\ D(s-1)}$     $S = \#$ POINTS    QUAD $=$ LAGRANGE POINTS

## Lemma

$\mathcal{L}^2$ operator of ADER defined by Gauss–Lobatto or Gauss–Legendre points and quadrature (they coincide) with $s = M + 1$ stages satisfies $\underbrace{C(s-1)}$ and $\underbrace{D(s-1)}$.

## Proof (1/4).

- Interpolation with $\phi^j$ is exact for polynomials of degree $s - 1$. ⟵ ✓
- The quadrature is exact for polynomials of degree $2s - 3$.
  Recall that $\underline{\underline{A}} = \cancel{\underline{\underline{M}}}$. Condition $C(s-1)$ reads    $\underline{\underline{A}} = \underline{\underline{n}}^{-1}\underline{\underline{R}}$

$$\underline{\underline{A}}\,\underline{c}^{z-1} \overset{!}{=} \frac{1}{z}\underline{c}^z \iff \underline{\underline{R}}\,\underline{c}^{z-1} \overset{?}{=} \frac{1}{z}\underline{\underline{M}}\,\underline{c}^z \iff \underline{\mathcal{X}} := \underline{\underline{R}}\,\underline{c}^{z-1} - \frac{1}{z}\underline{\underline{M}}\,\underline{c}^z \overset{?}{=} \underline{0}, \qquad z = 1, \ldots, s-1.$$

- Recall $\underline{\mathcal{F}}_m = t^m$, $\underline{b}_m = \underline{w}_m$, $\underline{\underline{R}}_{i,j} = \delta_{i,j} w_i$ and the definition of $\underline{\underline{M}}_{i,j} = \phi_i(1)\,\phi_j(\cdot) - \underbrace{\int_0^1 \phi_i'\,\phi_j}_{\text{exact}}$

$$\mathcal{X}_m := w_m \underbrace{(t^m)^{z-1}}_{} - \frac{1}{z}\left(\phi^m(1)\underset{\uparrow}{\phi^j(1)}\underbrace{(t^j)^z} - \int_0^1 \frac{d}{d\xi}\phi^m(\xi)\underbrace{\phi^j(\xi)(t^j)^z}\,d\xi\right).$$

$$\sum_j \phi^j(\xi)\cdot(\xi^j)^z = \xi^z \in \mathbb{P}^{s-1} \checkmark$$

$\underline{C(s-1)} \quad \underline{D(s-1)}$

Now, the interpolation of $t^z$ with $z \leq s-1$ with basis functions $\phi^j$ is exact. Hence, we can substitute $\phi^j(\xi)(t^j)^z = \xi^z$ for all $z = 1, \ldots, s-1$, obtaining

$$\mathcal{X}_m = w_m(t^m)^{z-1} - \frac{1}{z}\left(\phi^m(1)1^z - \int_0^1 \frac{d}{d\xi}\phi^m(\xi)\xi^z d\xi\right).$$

$$\mathbb{P}^{s-2} \quad \mathbb{P}^{s-1} \to \left[\phi^\wedge \, \xi^z\right]_0^1$$

Using the exactness of the quadrature for polynomials of degree $2s-3$, both true for Gauss–Lobatto and Gauss–Legendre, we know that the previous integral is exactly computed as $\frac{d}{d\xi}\phi^m(\xi)$ is of degree at most $s-2$ and $\xi^z$ is at most $s-1$. So, we can use integration by parts and obtain

$$\mathcal{X}_m = w_m(t^m)^{z-1} - \frac{1}{z}\left(\phi^m(0)0 + \int_0^1 \phi^m(\xi)\frac{d}{d\xi}\xi^z d\xi\right) = w_m(t^m)^{z-1} - \int_0^1 \phi^m(\xi)\xi^{z-1}d\xi = 0$$

$$w^m \cdot (\xi^\wedge)^{z-1}$$

by the exactness of the quadrature rule and the definition of $w_m$. Note that the condition is sharp, since the interpolation is not anymore exact for $z = s$, hence $\underline{C(s)}$ is not satisfied.

$$C(s-1) \checkmark$$

## $C(s-1)\ D(s-1)$

To prove $D(s-1)$, we write explicitly the condition in matricial form, for all $z = 1, \ldots, s-1$

$$\underline{\underline{bc^{z-1}}} \, \underline{\underline{A}} = \frac{1}{z}\underline{b}(1-c^z) \iff \underline{bc^{z-1}}\,\underline{\underline{M}}^{-1}\,\underline{\underline{R}} = \frac{1}{z}\underline{b}(1-c^z) \iff \underline{bc^{z-1}} = \frac{1}{z}\underline{b}(1-c^z)\underline{\underline{R}}^{-1}\,\underline{\underline{M}}.$$

Note that $b^m = w_m$ and $\underline{\underline{R}}^m_r = w_m \delta^m_r$, so $\underline{b}(1-c^z)\underline{\underline{R}}^{-1} = (1-c^z)$. It is left to prove that

$$\mathcal{Y} := \underline{bc^{z-1}} - \frac{1}{z}(1-c^z)\underline{\underline{M}} = \underline{0}.$$

$$\mathcal{Y}_m = w_m(t^m)^{z-1} - \frac{1}{z}\sum_{j=1}^{s}\left(1-(t^j)^z\right)\left(\phi^j(1)\phi^m(1) - \int_0^1 \frac{d}{d\xi}\phi^j(\xi)\phi^m(\xi)d\xi\right).$$

$$1 - \xi^z \in \mathbb{P}^{s-1}$$

$$\frac{d}{d\xi}(1-\xi^t) \in \mathbb{P}^{s-2}$$

$\underline{C(s-1)}$ ✓ $\underline{D(s-1)}$ ✓

## Proof (4/4).

Let us observe that, since $z \leq s - 1$, the polynomial is exactly represented by the Lagrangian interpolation $t^z = \sum_{j=1}^{s} \phi(t)(t^m)^z$. Hence, using the exactness of the quadrature for polynomials of degree at most $2s - 3$, we have

$$\mathcal{Y}_m = w_m (t^m)^{z-1} - \frac{1}{z} \left(1 - (1)^z\right) \phi^m(1) + \frac{1}{z} \int_0^1 \frac{d}{d\xi} \left(1 - (\xi)^z\right) \phi^m(\xi) d\xi$$

$$= w_m (t^m)^{z-1} - \frac{1}{z} \int_0^1 z \xi^{z-1} \phi^m(\xi) d\xi = w_m (t^m)^{z-1} - w_m (t^m)^{z-1} = 0. \quad \checkmark$$

Hence, ADER-Legendre and ADER-Lobatto satisfy $D(s-1)$. Note that the condition is sharp, since the interpolation is not anymore exact for $z = s$, hence $D(s)$ is not satisfied.

### Remark (ADER-Legendre is no collocation method)

*From the proof of previous Lemma, we can observe that ADER-Legendre methods do not satisfy $\overline{C(s)}$, hence, the methods are not collocation methods and they do not coincide with Gauss–Legendre implicit RK methods.*

$$\text{ORDER } 2S \Rightarrow C(S) \qquad \Rightarrow GLG \text{ ORDER} < 2S$$

### Theorem

$\mathcal{L}^2$ *of ADER with Gauss–Legendre is of order* $2s - 1$.

### Proof.

ADER-Legendre with $s = M + 1$ stages satisfies $\overline{B(2s)}$ for the quadrature rule and, hence, it satisfies $\overline{B(2s - 1)}$. For previous Lemma it also satisfies $C(s - 1)$ and $D(s - 1)$. Hence, Butcher's (1964) Theorem ($p \leq \eta + \zeta + 1$ and $p \leq 2\eta + 2$) guarantees that the method is of order $2s - 1$, since it is satisfied with $p = 2s - 1$ and $\eta = \zeta = s - 1$.

$$B(p) \quad C(\eta) \quad D(\zeta)$$

$$2s - 1 \leq s - 1 + s - 1 + 1 = 2s - 1 \checkmark \qquad p = 2s - 1 \leq 2(s - 1) + 2 = 2s \checkmark$$

# ADER Gauss–Lobatto $\mathcal{L}^2$

## Theorem

$\mathcal{L}^2$ of ADER with Gauss-Lobatto is of order $2s - 2$.

## Proof.

The condition for $B(2s - 2)$ is satisfied as $(c, b)$ is the Gauss–Lobatto quadrature with order $2s - 2$. Previous Lemma guarantees that ADER-Lobatto satisfies $B(2s - 2)$, $C(s - 1)$ and $D(s - 1)$, so Butcher's (1964) Theorem ($p \leq \eta + \zeta + 1$ and $p \leq 2\eta + 2$) is satisfied for order $p = 2s - 2$ and $\eta = \zeta = s - 1$.

$$p = 2s-2 \leq s-1+s-1+1 = 2s-1 \checkmark$$

$$p = 2s-2 \leq 2(s-1)+2 = 2s \checkmark$$

# ADER Gauss–Lobatto $\mathcal{L}^2$

## Theorem

$\mathcal{L}^2$ of ADER with Gauss-Lobatto is LobattoIIIC.

The Lobatto IIIC method is defined using the condition

$$a_{i1} = b_1 \quad \text{for } i = 1, \ldots, s. \tag{23}$$

*(handwritten annotations: $w_1$ above $b_1$; $+(C(s-1))$ with checkmark above the equation; arrow pointing to $a_{i1}$)*

## Lemma

$\mathcal{L}^2$ of ADER with Gauss-Lobatto satisfies (23).

## Theorem (Chipman 1971)

Lobatto IIIC schemes (in particular RK $a_{ij}$) are uniquely determined by Gauss–Lobatto quadrature rule $(c, b)$, condition (23) and by $C(s-1)$.

**Lemma**

$\mathcal{L}^2$ of ADER with Gauss-Lobatto satisfies (23).

**Proof.**

$$a_{i1} = \sum_j (\underline{\underline{M}}^{-1})_{ij} \mathbb{R}_{j1} = b_1 = w_1 \Longleftrightarrow$$

$$\sum_{i,j} \underline{\underline{M}}_{ki} (\underline{\underline{M}}^{-1})_{ij} \mathbb{R}_{j1} = \sum_i \underline{\underline{M}}_{ki} w_1 \Longleftrightarrow$$

$$\delta_{k1} w_1 = \mathbb{R}_{k1} = \sum_i \underline{\underline{M}}_{ki} w_1$$

$$\sum_i \underline{\underline{M}}_{ki} w_1 = \phi^m(1) w_1 - \int_0^1 \frac{d}{dt} \phi^m(\xi) w_1 \, dt = w_1 \phi^m(0) = w_1 \delta_{m,1}.$$

$\square$

## Outline

*DeC ADER*

## Usages

- Hyperbolic PDEs as explicit iterative methods (ADER: Toro, Dumbser, Klingenberg, Boscheri; DeC: Abgrall, Ricchiuto)
- IMEX solvers for hyperbolic with stiff sources (ADER: Dumbser, Boscheri; DeC: Abgrall, Torlo)
- IMEX solvers for hyperbolic with viscosity (treated implicitly) as compressible Navier Stokes (DeC: Minion, Dumbser, Zeifang)

## IMEX

$\partial_t u = F(u) + S(u)$
$S(u)$ stiff to be treated implicitly

## Advantages

- Arbitrary high order
- Unique framework to have matching between implicit and explicit terms
- Easy to code
- Iterative solver automatically included

## Disadvantages

- Explicit solver: many many stages
- Implicit: many stages
- Explicit: not amazing stability property (wrt SSP RK e.g.)

$$y'(t) = -|y(t)|y(t),$$
$$y(0) = 1, \qquad (24)$$
$$t \in [0, 0.1].$$

Convergence curves for ADER and DeC, varying the approximation order and collocation of nodes for the subtimesteps for a scalar nonlinear ODE

## Lotka–Volterra



Figure: Numerical solution of the Lotka-Volterra system using ADER (top) and DeC (bottom) with Gauss-Lobatto nodes with timestep $\Delta T = 1$.

Figure: Convergence error for Burgers equations: Left ADER right DeC. Space

## Outline

# Reduce computational cost for explicit DeC

## Literature

- *Micalizzi, L., Torlo, D. A new efficient explicit Deferred Correction framework: analysis and applications to hyperbolic PDEs and adaptivity.* arxiv.org/abs/2210.02976
- *Micalizzi, L., Torlo, D., Boscheri, W. Efficient iterative arbitrary high order methods: an adaptive bridge between low and high order.* arxiv.org/abs/2212.07783

## Goal

Reduce computational costs of explicit DeC.

## DeC as RK for ODEs

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(p)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(p-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(p-1)}) \text{ with } p = 1, \ldots, P.$$

$$\boldsymbol{u}^{m,(p)} = \boldsymbol{u}^0 + \sum_{r=0}^{M} \theta_r^m F(t^r, \boldsymbol{u}^{r,(p-1)}), \qquad \forall m = 1, \ldots, M, \ p = 1, \ldots, P$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(p)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(p-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(p-1)}) \text{ with } p = 1, \ldots, P.$$

$$\boldsymbol{u}^{m,(p)} = \boldsymbol{u}^0 + \sum_{r=0}^{M} \theta_r^m F(t^r, \boldsymbol{u}^{r,(p-1)}), \qquad \forall m = 1, \ldots, M, \ p = 1, \ldots, P$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(p)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(p-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(p-1)}) \text{ with } p = 1, \ldots, P.$$

$$\boldsymbol{u}^{m,(p)} = \boldsymbol{u}^0 + \sum_{r=0}^{M} \theta_r^m F(t^r, \boldsymbol{u}^{r,(p-1)}), \qquad \forall m = 1, \ldots, M, \; p = 1, \ldots, P$$

$$\mathcal{L}^1(\underline{\boldsymbol{u}}^{(p)}) = \mathcal{L}^1(\underline{\boldsymbol{u}}^{(p-1)}) - \mathcal{L}^2(\underline{\boldsymbol{u}}^{(p-1)}) \text{ with } p = 1, \ldots, P.$$

$$\boldsymbol{u}^{m,(p)} = \boldsymbol{u}^0 + \sum_{r=0}^{M} \theta_r^m F(t^r, \boldsymbol{u}^{r,(p-1)}), \qquad \forall m = 1, \ldots, M, \ p = 1, \ldots, P$$



| $\underline{c}$ | $\boldsymbol{u}^0$ | $\boldsymbol{u}^{(1)}$ | $\boldsymbol{u}^{(2)}$ | $\boldsymbol{u}^{(3)}$ | $\cdots$ | $\boldsymbol{u}^{(M-1)}$ | $\boldsymbol{u}^{(M)}$ | A |
|---|---|---|---|---|---|---|---|---|
| 0 | 0 | | | | | | | $\boldsymbol{u}^0$ |
| $\underline{\beta}_{1:}$ | $\underline{\beta}_{1:}$ | $\underline{\underline{0}}$ | | | | | | $\boldsymbol{u}^{(1)}$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | $\boldsymbol{u}^{(2)}$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\underline{\underline{0}}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | | | | $\boldsymbol{u}^{(3)}$ |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\ddots$ | $\ddots$ | | | $\vdots$ |
| | $\vdots$ | $\vdots$ | | | $\ddots$ | $\ddots$ | | $\vdots$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | $\boldsymbol{u}^{(M)}$ |
| $\underline{b}$ | $\Theta_{M,0}$ | $\underline{0}$ | $\cdots$ | $\cdots$ | $\cdots$ | $\underline{0}$ | $\Theta_{M,1:}$ | $\boldsymbol{u}^{M,(M+1)}$ |

# Costs

**Large costs!**

## Large costs!

Equispaced

| P | M | DeC |
|---|---|-----|
| 2 | 1 | 2 |
| 3 | 2 | 5 |
| 4 | 3 | 10 |
| 5 | 4 | 17 |
| 6 | 5 | 26 |
| 7 | 6 | 37 |
| 8 | 7 | 50 |
| 9 | 8 | 65 |
| 10 | 9 | 82 |

Gauss–Lobatto

| P | M | DeC |
|---|---|-----|
| 2 | 1 | 2 |
| 3 | 2 | 5 |
| 4 | 2 | 7 |
| 5 | 3 | 13 |
| 6 | 3 | 16 |
| 7 | 4 | 25 |
| 8 | 4 | 29 |
| 9 | 5 | 41 |
| 10 | 5 | 46 |

- DeC $S = M \cdot (P-1) + 1$
  - DeC equi $S = (P-1)^2 + 1$
  - DeC GLB $S = \left\lceil \frac{P}{2} \right\rceil (P-1) + 1$

## Large costs!

- DeC $S = M \cdot (P - 1) + 1$
  - DeC equi $S = (P - 1)^2 + 1$
  - DeC GLB $S = \left\lceil \frac{P}{2} \right\rceil (P - 1) + 1$

**Equispaced**

| $P$ | $M$ | DeC |
|-----|-----|-----|
| 2 | 1 | 2 |
| 3 | 2 | 5 |
| 4 | 3 | 10 |
| 5 | 4 | 17 |
| 6 | 5 | 26 |
| 7 | 6 | 37 |
| 8 | 7 | 50 |
| 9 | 8 | 65 |
| 10 | 9 | 82 |

**Gauss–Lobatto**

| $P$ | $M$ | DeC |
|-----|-----|-----|
| 2 | 1 | 2 |
| 3 | 2 | 5 |
| 4 | 2 | 7 |
| 5 | 3 | 13 |
| 6 | 3 | 16 |
| 7 | 4 | 25 |
| 8 | 4 | 29 |
| 9 | 5 | 41 |
| 10 | 5 | 46 |

**How can we save computational time?**

## Outline

## Idea for reduction of stages

Order $\quad O(\Delta t^1) \quad O(\Delta t^2) \quad O(\Delta t^3) \quad O(\Delta t^4) \quad O(\Delta t^5) \quad O(\Delta t^6)$

$t_n = t^0$

$\mathbf{u}^{0,(0)} \quad \mathbf{u}^{0,(1)} \quad \mathbf{u}^{0,(2)} \quad \mathbf{u}^{0,(3)} \quad \mathbf{u}^{0,(4)} \quad \mathbf{u}^{0,(5)}$

$t^1$

$\mathbf{u}^{1,(4)} \quad \mathbf{u}^{1,(5)}$

$t^2$

$\mathbf{u}^{1,(2)} \quad \mathbf{u}^{1,(3)} \quad \mathbf{u}^{2,(4)} \quad \mathbf{u}^{2,(5)}$

$t^3$

$\mathbf{u}^{2,(3)} \quad \mathbf{u}^{3,(4)} \quad \mathbf{u}^{3,(5)}$

$t_{n+1} = t^M = t^4$

$\mathbf{u}^{1,(0)} \quad \mathbf{u}^{1,(1)} \quad \mathbf{u}^{2,(2)} \quad \mathbf{u}^{3,(3)} \quad \mathbf{u}^{4,(4)} \quad \mathbf{u}^{4,(5)}$

Iteration $\quad 0 \qquad 1 \qquad 2 \qquad 3 \qquad 4 \qquad P=5$

**DeC**

$$\underline{\boldsymbol{u}}^{(p)} = \underline{\boldsymbol{u}}^0 + \Delta t \Theta F(\underline{\boldsymbol{u}}^{(p-1)})$$

## How to communicate between iterations?



**DeC**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta F(\underline{u}^{(p-1)})$$

**DeCu**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta^{(p)} F(H^{(p)} \underline{u}^{(p-1)})$$

$$H_{ij}^{(p)} = \phi_j^{(p-1)}(t^{i,(p)})$$

**DeC**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta F(\underline{u}^{(p-1)})$$

**DeCu**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta^{(p)} F(H^{(p)} \underline{u}^{(p-1)})$$

**DeCdu**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta^{(p)} H^{(p)} F(\underline{u}^{(p-1)})$$

$$H_{ij}^{(p)} = \phi_j^{(p-1)}(t^{i,(p)})$$

**DeC**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta F(\underline{u}^{(p-1)})$$

**DeCu**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta^{(p)} F(H^{(p)} \underline{u}^{(p-1)})$$
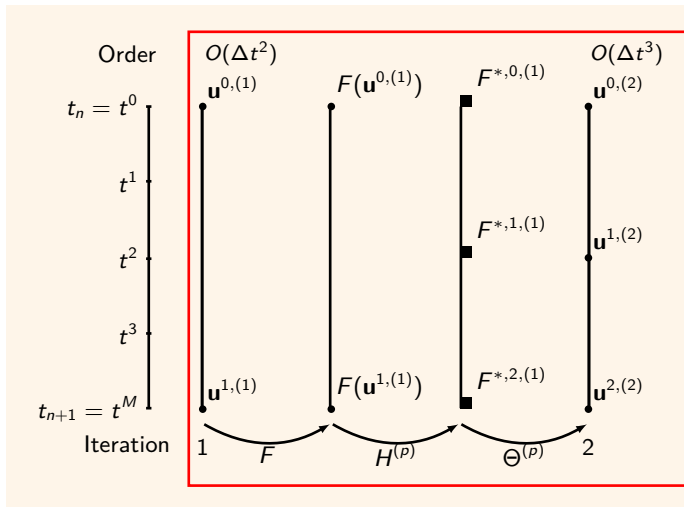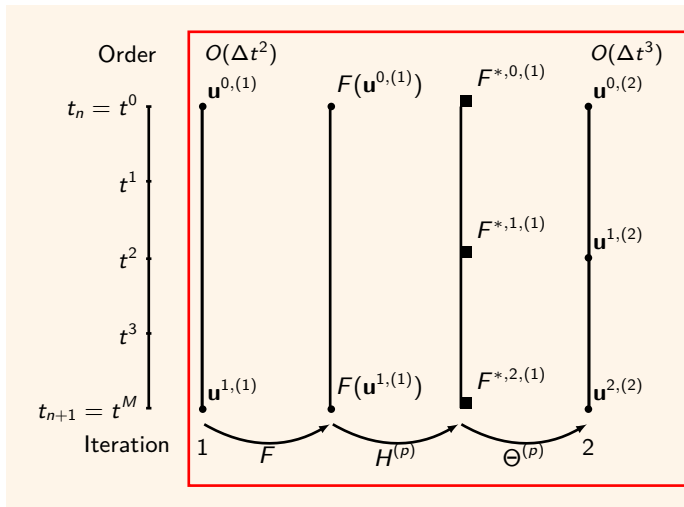
$$\underline{u}^{*(p)} = \underline{u}^0 + \Delta t H^{(p)} \Theta^{*(p-1)} F(\underline{u}^{*(p-1)})$$

**DeCdu**

$$\underline{u}^{(p)} = \underline{u}^0 + \Delta t \Theta^{(p)} H^{(p)} F(\underline{u}^{(p-1)})$$

$$H_{ij}^{(p)} = \phi_j^{(p-1)}(t^{i,(p)})$$

## DeC $\quad S = M \cdot (P-1) + 1$

| $\underline{c}$ | $\mathbf{u}^0$ | $\mathbf{u}^{(1)}$ | $\mathbf{u}^{(2)}$ | $\mathbf{u}^{(3)}$ | $\cdots$ | $\mathbf{u}^{(M-1)}$ | $\mathbf{u}^{(M)}$ | A | dim |
|---|---|---|---|---|---|---|---|---|---|
| $0$ | $0$ | | | | | | | $\mathbf{u}^0$ | $1$ |
| $\underline{\beta}_{1:}$ | $\underline{\beta}_{1:}$ | $\underline{\underline{0}}$ | | | | | | $\mathbf{u}^{(1)}$ | $M$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | $\mathbf{u}^{(2)}$ | $M$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\underline{\underline{0}}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | | | | $\mathbf{u}^{(3)}$ | $M$ |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\ddots$ | $\ddots$ | | | $\vdots$ | $M$ |
| | $\vdots$ | $\vdots$ | | | $\ddots$ | $\ddots$ | | $\vdots$ | $M$ |
| $\underline{\beta}_{1:}$ | $\Theta_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $\Theta_{1:,1:}$ | $\underline{\underline{0}}$ | $\mathbf{u}^{(M)}$ | $M$ |
| $\underline{b}$ | $\Theta_{M,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $\Theta_{M,1:}$ | $\mathbf{u}^{M,(M+1)}$ | |

**DeCu** $\quad S = M \cdot (P - 1) + 1 - \frac{(M-1)(M-2)}{2}$

| $\underline{c}$ | $\mathbf{u}^0$ | $\mathbf{u}^{*(1)}$ | $\mathbf{u}^{*(2)}$ | $\mathbf{u}^{*(3)}$ | $\cdots$ | $\mathbf{u}^{*(M-2)}$ | $\mathbf{u}^{*(M-1)}$ | $\mathbf{u}^{(M)}$ | A | dim |
|---|---|---|---|---|---|---|---|---|---|---|
| $0$ | $0$ | | | | | | | | $\mathbf{u}^0$ | $1$ |
| $\underline{\beta}^{(2)}_{1:}$ | $\underline{\beta}^{(2)}_{1:}$ | $\underline{\underline{0}}$ | | | | | | | $\mathbf{u}^{*(1)}$ | $2$ |
| $\underline{\beta}^{(3)}_{1:}$ | $W^{(2)}_{1:,0}$ | $W^{(2)}_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | | $\mathbf{u}^{*(2)}$ | $3$ |
| $\underline{\beta}^{(4)}_{1:}$ | $W^{(3)}_{1:,0}$ | $\underline{\underline{0}}$ | $W^{(3)}_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | $\mathbf{u}^{*(3)}$ | $4$ |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\ddots$ | $\ddots$ | | | | $\vdots$ | $\vdots$ |
| | | | | | $\ddots$ | $\ddots$ | | | $\vdots$ | $\vdots$ |
| $\underline{\beta}^{(M)}_{1:}$ | $W^{(M-1)}_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $W^{(M-1)}_{1:,1:}$ | $\underline{\underline{0}}$ | $\underline{\underline{0}}$ | $\mathbf{u}^{*(M-1)}$ | $M$ |
| $\underline{\beta}^{(M)}_{1:}$ | $W^{(M)}_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $W^{(M)}_{1:,1:}$ | $\underline{\underline{0}}$ | | $\mathbf{u}^{(M)}$ | $M$ |
| $\underline{b}$ | $W^{(M+1)}_{M,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $W^{(M+1)}_{M,1:}$ | $\mathbf{u}^{M,(M+1)}$ | | |

$$W^{(p)} := \begin{cases} H^{(p)}\Theta^{(p)} \in \mathbb{R}^{(p+2)\times(p+1)}, & \text{if } p = 2, \ldots, M-1, \\ \Theta^{(M)} \in \mathbb{R}^{(M+1)\times(M+1)}, & \text{if } p \geq M. \end{cases}$$

# Efficient DeC into RK framework

**DeCdu** $\qquad S = M \cdot (P-1) + 1 - \dfrac{M(M-1)}{2}$

| $\underline{c}$ | $\mathbf{u}^0$ | $\mathbf{u}^{(1)}$ | $\mathbf{u}^{(2)}$ | $\mathbf{u}^{(3)}$ | $\cdots$ | $\mathbf{u}^{(M-2)}$ | $\mathbf{u}^{(M-1)}$ | $\mathbf{u}^{(M)}$ | A | dim |
|---|---|---|---|---|---|---|---|---|---|---|
| $0$ | $0$ | | | | | | | | $\mathbf{u}^0$ | 1 |
| $\underline{\beta}^{(1)}_{1:}$ | $\underline{\beta}^{(1)}_{1:}$ | $\underline{\underline{0}}$ | | | | | | | $\mathbf{u}^{(1)}$ | 1 |
| $\underline{\beta}^{(2)}_{1:}$ | $Z^{(2)}_{1:,0}$ | $Z^{(2)}_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | | $\mathbf{u}^{(2)}$ | 2 |
| $\underline{\beta}^{(3)}_{1:}$ | $Z^{(3)}_{1:,0}$ | $\underline{\underline{0}}$ | $Z^{(3)}_{1:,1:}$ | $\underline{\underline{0}}$ | | | | | $\mathbf{u}^{(3)}$ | 3 |
| $\vdots$ | $\vdots$ | $\vdots$ | | $\ddots$ | $\ddots$ | | | | $\vdots$ | $\vdots$ |
| | | $\vdots$ | | | $\ddots$ | $\ddots$ | | | $\vdots$ | $\vdots$ |
| $\underline{\beta}^{(M-1)}_{1:}$ | $Z^{(M-1)}_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $Z^{(M-1)}_{1:,1:}$ | $\underline{\underline{0}}$ | $\underline{\underline{0}}$ | $\mathbf{u}^{(M-1)}$ | $M-1$ |
| $\underline{\beta}^{(M)}_{1:}$ | $Z^{(M)}_{1:,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $Z^{(M)}_{1:,1:}$ | $\underline{\underline{0}}$ | $\mathbf{u}^{(M)}$ | $M$ |
| $\underline{b}$ | $Z^{(M+1)}_{M,0}$ | $\underline{\underline{0}}$ | $\cdots$ | $\cdots$ | $\cdots$ | $\cdots$ | $\underline{\underline{0}}$ | $Z^{(M+1)}_{M,1:}$ | $\mathbf{u}^{M,(M+1)}$ | |

$$Z^{(p)} := \begin{cases} \Theta^{(p)} H^{(p-1)} \in \mathbb{R}^{(p+1)\times p}, & \text{if } p = 1,\dots,M, \\ \Theta^{(M)} \in \mathbb{R}^{(M+1)\times(M+1)}, & \text{if } p > M. \end{cases}$$

## Equispaced

| P | M | DeC | DeCu | DeCdu |
|---|---|-----|------|-------|
| 2 | 1 | 2 | 2 | 2 |
| 3 | 2 | 5 | 5 | 4 |
| 4 | 3 | 10 | 9 | 7 |
| 5 | 4 | 17 | 14 | 11 |
| 6 | 5 | 26 | 20 | 16 |
| 7 | 6 | 37 | 27 | 22 |
| 8 | 7 | 50 | 35 | 29 |
| 9 | 8 | 65 | 44 | 37 |
| 10 | 9 | 82 | 54 | 46 |
| 11 | 10 | 101 | 65 | 56 |
| 12 | 11 | 122 | 77 | 67 |
| 13 | 12 | 145 | 90 | 79 |

## Gauss–Lobatto

| P | M | DeC | DeCu | DeCdu |
|---|---|-----|------|-------|
| 2 | 1 | 2 | 2 | 2 |
| 3 | 2 | 5 | 5 | 4 |
| 4 | 2 | 7 | 7 | 6 |
| 5 | 3 | 13 | 12 | 10 |
| 6 | 3 | 16 | 15 | 13 |
| 7 | 4 | 25 | 22 | 19 |
| 8 | 4 | 29 | 26 | 23 |
| 9 | 5 | 41 | 35 | 31 |
| 10 | 5 | 46 | 40 | 36 |
| 11 | 6 | 61 | 51 | 46 |
| 12 | 6 | 67 | 57 | 52 |
| 13 | 7 | 85 | 70 | 64 |

**DeC, DeCu, DeCdu**

## DeC-DeCu-DeCdu

The **stability function** of DeC, DeCu, DeCdu of order $P$ for any nodes distribution is

$$R(z) = 1 + z + \frac{z^2}{2!} + \cdots + \frac{z^P}{P!}.$$

## Efficient DeC

- Code DeCu or DeCdu
- Check order of accuracy
- Write a code to obtain its RK matrix
- Check the stability function with nodepy
- Compare computational costs with original DeC