# High order IMEX deferred correction residual distribution schemes for stiff kinetic problems

Davide Torlo

Team Cardamom
INRIA Bordeaux – Sud-Ouest

26th November 2020

joint work with Rèmi Abgrall and Mario Ricchiuto

# My research

## Education

- PostDoc: INRIA, prof. Mario Ricchiuto
- PhD: University of Zurich, prof. Rémi Abgrall
- Master: SISSA Trieste, prof. Gianluigi Rozza
- Bachelor: Università di Milano–Bicocca

## Research

- Model order reduction (advection dominated problems)
- High order methods for hyperbolic problems (kinetic problems)
- High order methods for positive ODEs
- Structure preserving methods

# My research

## Education

- PostDoc: INRIA, prof. Mario Ricchiuto
- PhD: University of Zurich, prof. Rémi Abgrall
- Master: SISSA Trieste, prof. Gianluigi Rozza
- Bachelor: Università di Milano–Bicocca

## Research

- Model order reduction (advection dominated problems)
- High order methods for hyperbolic problems (**kinetic problems**)
- High order methods for positive ODEs
- Structure preserving methods

# Outline

## Motivation: relaxed systems

What we want to solve is an hyperbolic relaxation system:

$$\partial_t u + \nabla_x \cdot A(u) = \frac{S(u)}{\varepsilon} \text{ or}$$
$$\partial_t u + H(u)\nabla_x u = \frac{S(u)}{\varepsilon}$$

(1)

Applications:

- Jin–Xin system
- Kinetic models
- Multiphase flows
- Viscoelasticity problems

What we want to solve is an hyperbolic relaxation system:

$$\partial_t u + \nabla_x \cdot A(u) = \frac{S(u)}{\varepsilon} \text{ or}$$
$$\partial_t u + H(u)\nabla_x u = \frac{S(u)}{\varepsilon} \tag{1}$$

Applications:

- Jin–Xin system
- Kinetic models
- Multiphase flows
- Viscoelasticity problems

## Goal

A scheme that is

- Asymptotic preserving:

$$
\begin{array}{ccc}
\mathcal{F}_\Delta^\varepsilon & \xrightarrow{\varepsilon \to 0} & \mathcal{F}_\Delta^0 \\
\Delta \to 0 \Big\downarrow & & \Big\downarrow \Delta \to 0 \\
\mathcal{F}^\varepsilon & \xrightarrow{\varepsilon \to 0} & \mathcal{F}^0
\end{array}
$$

- High order in space and time
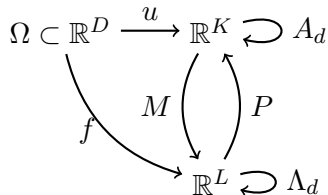- Computationally explicit (as much as possible, no mass matrix)

# Outline

## Kinetic Models

Kinetic relaxation models by D. Aregba-Driollet and R. Natalini[1].
Hyperbolic limit equation is

$$u_t + \sum_{d=1}^{D} \partial_{x_d} A_d(u) = 0, \quad u : \Omega \to \mathbb{R}^K.$$

$$\Omega \subset \mathbb{R}^D \xrightarrow{u} \mathbb{R}^K \circlearrowright A_d$$

$$f \searrow \quad M \updownarrow \quad \updownarrow P$$

$$\mathbb{R}^L \circlearrowright \Lambda_d$$

Relaxation system

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right), \quad f^\varepsilon : \Omega \to \mathbb{R}^L$$

$$Pf^\varepsilon \to u, \quad P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

---

[1] D. Aregba-Driollet and R. Natalini. Discrete kinetic schemes for multidimensional systems of conservation laws. SIAM J. Numer. Anal., 37(6):1973–2004, 2000.

# Chapman–Enskog

Relaxation system

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right),$$

$$P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

Define $u^\varepsilon = Pf^\varepsilon$, $v_d^\varepsilon = P\Lambda_d f^\varepsilon$

$$\begin{cases} \partial_t u^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} v_j^\varepsilon = 0 \\ \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j}(P\Lambda_j \Lambda_d f^\varepsilon) = \frac{1}{\varepsilon}(A_d(u^\varepsilon) - v_d^\varepsilon), \end{cases}$$

## Chapman–Enskog

Relaxation system

Define $u^\varepsilon = Pf^\varepsilon$, $v_d^\varepsilon = P\Lambda_d f^\varepsilon$

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon}\left(M(Pf^\varepsilon) - f^\varepsilon\right),$$

$$P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

$$\begin{cases} \partial_t u^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} v_j^\varepsilon = 0 \\ \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j}(P\Lambda_j\Lambda_d f^\varepsilon) = \frac{1}{\varepsilon}(A_d(u^\varepsilon) - v_d^\varepsilon), \end{cases}$$

$$v_d^\varepsilon = A_d(u^\varepsilon) - \varepsilon\left(\partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j}(P\Lambda_d\Lambda_j M(u^\varepsilon))\right) + \mathcal{O}(\varepsilon^2),$$

$$\partial_t u^\varepsilon + \sum_{d=1}^{D} \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^{D} \partial_{x_d}\left(\sum_{j=1}^{D} B_{dj}(u^\varepsilon)\partial_{x_j} u^\varepsilon\right) + \mathcal{O}(\varepsilon^2)$$

with $B_{dj}(u) := P\Lambda_d\Lambda_j M'(u) - A_d'(u)A_j'(u) \in \mathbb{R}^{S\times S}$, $\forall d, j = 1, \ldots, D$.

$$\partial_t u^\varepsilon + \sum_{d=1}^{D} \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^{D} \partial_{x_d} \left( \sum_{j=1}^{D} B_{dj}(u^\varepsilon) \partial_{x_j} u^\varepsilon \right) + \mathcal{O}(\varepsilon^2).$$

Right hand side must be diffusive.
Whitham's subcharacteristic condition[2] becomes

$$B_{jd} := P\Lambda_d \Lambda_j M'(u) - A'_d(u)A'_j(u), \qquad \sum_{j,d=1}^{D} (B_{dj}\xi_j, \xi_d) \geq 0.$$

[2] D. Aregba-Driollet and R. Natalini. Discrete kinetic schemes for multidimensional systems of conservation laws. SIAM J. Numer. Anal., 37(6):1973–2004, 2000.

## Kinetic model

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right), \qquad P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

We have to find $M, P, \Lambda$ that respect previous conditions.
$L = N \times K$ with $P = (I_K, \dots, I_K)$ $N$ blocks of identity matrices in $\mathbb{R}^K$.
$f_n \in \mathbb{R}^K$ with $n = 1, \dots, N$

$$\Lambda_d = diag(\Lambda_1^{(d)}, \dots, \Lambda_N^{(d)}) \qquad \Lambda_n^{(d)} = \lambda_n^{(d)} I_K, \quad \text{for } \lambda_n^{(d)} \in \mathbb{R}.$$

With this formalism we can rewrite (43) as

$$\begin{cases} \partial_t f_n^\varepsilon + \sum_{d=1}^{D} \Lambda_n^{(d)} \partial_{x_d} f_n^\varepsilon = \frac{1}{\varepsilon} \left( M_n(u^\varepsilon) - f_n^\varepsilon \right), & \forall n = 1, \dots, N \\ u^\varepsilon = \sum_{n=1}^{N} f_n^\varepsilon \end{cases} \tag{2}$$

## Kinetic model – DRM

Let us present the *diagonal relaxation method (DRM)*. Here $N = D + 1$. Then we have to define maxwellians $M_n$ and matrices $\Lambda_j^{(d)}$. Take $\lambda > 0$ and

$$
\Lambda_j^{(d)} = \begin{cases} -\lambda I_K & j = d \\ \lambda I_K & j = D + 1 \\ 0 & \text{else} \end{cases} .
$$

The Maxwellians can be defined as follows:

$$
\begin{cases} M_{D+1}(u) = \left( u + \frac{1}{\lambda} \sum_{d=1}^{D} A_d(u) \right) / (D+1) \\ M_j(u) = -\frac{1}{\lambda} A_j(u) + M_{D+1}(u) \end{cases}
$$

Important: we have to choose $\lambda$ according to Whitham's subcharacteristic condition.

## Example of DMR model

$u : \Omega \subset \mathbb{R} \to \mathbb{R}, \quad D = 1, N = 2, \quad f : \mathbb{R} \to \mathbb{R}^2$

Limit equation

$$u_t + a(u)_x = 0 \tag{3}$$

$$\Lambda = \begin{pmatrix} -\lambda & 0 \\ 0 & \lambda \end{pmatrix}, \quad M(u) = \begin{pmatrix} \frac{u}{2} - \frac{a(u)}{2\lambda} \\ \frac{u}{2} + \frac{a(u)}{2\lambda} \end{pmatrix}, \quad Pf = f_1 + f_2 \tag{4}$$

Kinetic model is

$$\begin{cases} \partial_t f_1 - \lambda \partial_x f_1 = \frac{1}{\epsilon} \left( \frac{f_1 + f_2}{2} - \frac{a(f_1 + f_2)}{2\lambda} - f_1 \right) \\ \partial_t f_2 + \lambda \partial_x f_2 = \frac{1}{\epsilon} \left( \frac{f_1 + f_2}{2} + \frac{a(f_1 + f_2)}{2\lambda} - f_2 \right) \end{cases} \tag{5}$$

# Outline

# Residual Distribution

- High order
- Easy to code
- FE based
- Compact stencil
- No need of Riemann solver
- No need of conservative variables
- Can recast some other FV, FE schemes[3]

$$\partial_t f + \nabla_x \cdot A(f) = S(f)$$

$$V_h = \{f \in L^2(\Omega_h, \mathbb{R}^D) \cap \mathcal{C}^0(\Omega_h), \ f|_K \in \mathbb{P}^k, \ \forall K \in \Omega_h\}.$$

---

[3]R. Abgrall. Some remarks about conservation for residual distribution schemes. Computational Methods in Applied Mathematics, 2018. DOI: https://doi.org/10.1515/cmam-2017-0056.

Figure: Defining total residual, nodal residuals and building the RD scheme

Figure: Defining total residual, nodal residuals and building the RD scheme
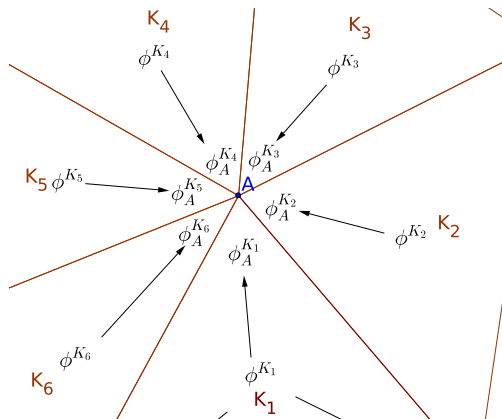
Figure: Defining total residual, nodal residuals and building the RD scheme

# Residual Distribution - Spatial Discretization

1. Define $\forall K \in \Omega_h$ a fluctuation term (total residual) $\phi^K = \int_K \nabla \cdot A(f) - S(f) dx$

2. Define a nodal residual $\phi_\sigma^K \ \forall \sigma \in K$ :

$$\phi^K = \sum_{\sigma \in K} \phi_\sigma^K, \quad \forall K \in \Omega_h. \tag{6}$$

3. The resulting scheme is

$$\partial_t f_\sigma + \sum_{K | \sigma \in K} \phi_\sigma^K = 0, \quad \forall \sigma \in D_h. \tag{7}$$

Remark: the definition of the nodal residuals leads to the scheme!
We use as Galerkin, Rusanov, PSI limiter, jump stabilization.

## Residual Distribution – Examples

How to split into $\phi_\sigma^K \Rightarrow$ choice of the scheme. For example, we can rewrite SUPG in this way:

$$\phi_\sigma^K(f) = \int_K \varphi_\sigma(\nabla \cdot A(f) - S(f))dx+ \tag{8}$$

$$+h_K \int_K (\partial_f A(f) \cdot \nabla \varphi_\sigma) \tau (\nabla \cdot A(f) - S(f)). \tag{9}$$

Furthermore, we can write the Galerkin FEM scheme with jump stabilization[4]:

$$\phi_\sigma^K = \int_K \varphi_\sigma(\nabla \cdot A(f) - S(f))dx + \sum_{e|\text{edge of } K} \theta h_e^2 \int_e [\nabla f] \cdot [\nabla \varphi_\sigma]d\Gamma, \tag{10}$$

---

[4] E. Burman and P. Hansbo. Comp. Meth. in Appl. Mech. and Eng., 193(15):1437 – 1453, 2004.

# Outline

# IMEX discretization - Kinetic model

Stiff source term $\Rightarrow$ oscillations when $\varepsilon \ll \Delta t$

$\Delta t \approx \varepsilon$ not feasible

IMEX approach: IMplicit for source term, EXplicit for advection term

$$\frac{f^{n+1,\varepsilon} - f^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^{n,\varepsilon} = \frac{1}{\varepsilon} \left( M(Pf^{n+1,\varepsilon}) - f^{n+1,\varepsilon} \right) \tag{11}$$

$$f^{0,\varepsilon}(x) = f_0^{\varepsilon}(x)$$

How to treat non-linear implicit functions?

Recall: $PM(u) = u$ and $Pf^{\varepsilon} = u^{\varepsilon}$, so

$$\frac{u^{n+1,\varepsilon} - u^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^{D} P\Lambda_d \partial_{x_d} f^{n,\varepsilon} = 0. \tag{12}$$

Find $u^{n+1,\varepsilon} = Pf^{n+1,\varepsilon}$ and substitute it in (11).

IMEX formulation = $\mathcal{L}^1$ (first order accurate).

## IMEX discretization - Kinetic model

Stiff source term $\Rightarrow$ oscillations when $\varepsilon \ll \Delta t$

$\Delta t \approx \varepsilon$ not feasible

IMEX approach: IMplicit for source term, EXplicit for advection term

$$\frac{f^{n+1,\varepsilon} - f^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^{n,\varepsilon} = \frac{1}{\varepsilon} \left( M(Pf^{n+1,\varepsilon}) - f^{n+1,\varepsilon} \right)$$

$$f^{0,\varepsilon}(x) = f_0^{\varepsilon}(x)$$

(11)

How to treat non-linear implicit functions?

Recall: $PM(u) = u$ and $Pf^{\varepsilon} = u^{\varepsilon}$, so

$$\frac{u^{n+1,\varepsilon} - u^{n,\varepsilon}}{\Delta t} + \sum_{d=1}^{D} P\Lambda_d \partial_{x_d} f^{n,\varepsilon} = 0.$$

(12)

Find $u^{n+1,\varepsilon} = Pf^{n+1,\varepsilon}$ and substitute it in (11).

IMEX formulation = $\mathcal{L}^1$ (first order accurate).

# IMEX is asymptotic preserving

To prove AP: induction.

$$\begin{array}{ccc} \mathcal{F}_\Delta^\varepsilon & \xrightarrow{\varepsilon \to 0} & \mathcal{F}_\Delta^0 \\ {\scriptstyle \Delta \to 0} \Big\downarrow & & \Big\downarrow {\scriptstyle \Delta \to 0} \\ \mathcal{F}^\varepsilon & \xrightarrow{\varepsilon \to 0} & \mathcal{F}^0 \end{array}$$

### Induction Hypothesis

$$\frac{u^{n+1} - u^n}{\Delta t} + \sum_{d=1}^{D} \partial_{x_d} A_d(u^n) + \mathcal{O}(\varepsilon) + \mathcal{O}(\Delta) = 0 \tag{13}$$

$$\frac{f^{n+1} - f^n}{\Delta t} + \sum_{d=1}^{D} \partial_{x_d} \Lambda_d f^n - \frac{M(u^{n+1}) - f^{n+1}}{\varepsilon} + \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta) = 0 \tag{14}$$

Given that the space discretization is consistent with the model.

# Outline

# DeC high order time discretization: $\mathcal{L}^2$
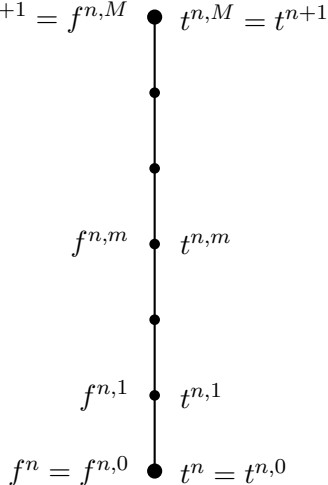
High order in time: we discretize our variable on $[t^n, t^{n+1}]$ in $M$ substeps ($f_\sigma^{n,m}$).

$$f^{n+1} = f^{n,M} \quad \bullet \quad t^{n,M} = t^{n+1}$$

Thanks to Picard–Lindelöf theorem, we can rewrite

$$f_\sigma^{n,m} = f_\sigma^{n,0} + \int_{t^n}^{t^{n,m}} \nabla \cdot A(f(x,s)) - S(f(x,s))ds$$
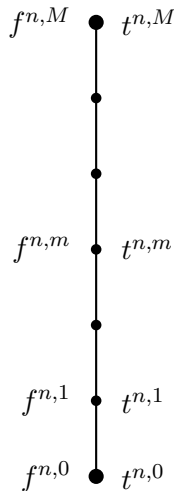
and if we want to reach order $r + 1$ we need $M = r$.

$$f^{n,m} \quad \bullet \quad t^{n,m}$$

$$f^{n,1} \quad \bullet \quad t^{n,1}$$

$$f^n = f^{n,0} \quad \bullet \quad t^n = t^{n,0}$$

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

$$
\mathcal{L}_\sigma^2(f^{n,0}, \ldots, f^{n,M}) =
$$

$$
= \begin{pmatrix} \sum_{K \ni \sigma} \left( \int\limits_K \varphi_\sigma(f^{n,M}(x) - f^{n,0}(x))dx + \Delta t \sum_{r=0}^M \theta_r^M \phi_\sigma^K(f^{n,r}) \right) \\ \vdots \\ \sum_{K \ni \sigma} \left( \int\limits_K \varphi_\sigma(f^{n,1}(x) - f^{n,0}(x))dx + \Delta t \sum_{r=0}^M \theta_r^1 \phi_\sigma^K(f^{n,r}) \right) \end{pmatrix}
$$

which is a fully implicit system of $M$ equations with $M$ unknowns (times #DoFs).

$f^{n,M} \bullet\ t^{n,M}$

$\bullet$

$\bullet$

$f^{n,m} \bullet\ t^{n,m}$

$\bullet$

$f^{n,1} \bullet\ t^{n,1}$

$f^{n,0} \bullet\ t^{n,0}$

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

$$\mathcal{L}_\sigma^1(f^{n,0}, \ldots, f^{n,M}) =$$

$$= \begin{pmatrix} \sum_{K \ni \sigma} \left( (f_\sigma^{n,M} - f_\sigma^{n,0}) \int\limits_K \varphi_\sigma dx + \Delta t \beta^M \phi_\sigma^K(f^{n,0}, f^{n,M}) \right) \\ \vdots \\ \sum_{K \ni \sigma} \left( (f_\sigma^{n,1} - f_\sigma^{n,0}) \int\limits_K \varphi_\sigma dx + \Delta t \beta^1 \phi_\sigma^K(f^{n,0}, f^{n,1}) \right) \end{pmatrix}$$

- IMEX discretization
- mass lumping on implicit terms (time derivative and source term)
- easy to be solved (explicit or small implicit systems)

$f^{n,M} \quad t^{n,M}$

$f^{n,m} \quad t^{n,m}$

$f^{n,1} \quad t^{n,1}$

$f^{n,0} \quad t^{n,0}$

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

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

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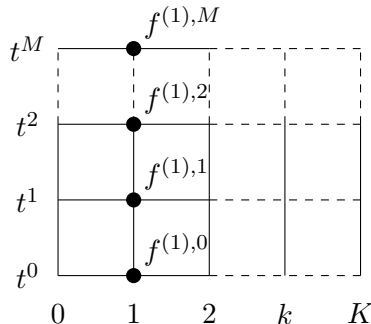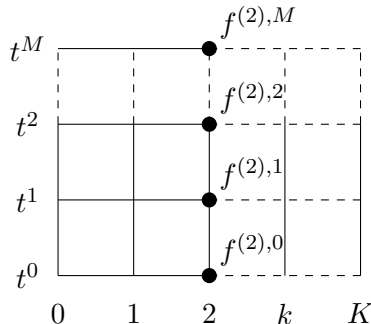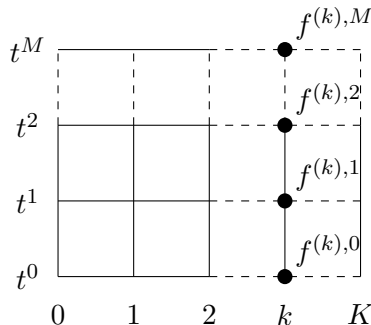$$f^{m,(0)} := f(t^n), \quad m = 1, \ldots, M$$
$$\mathcal{L}^1(f^{(k)}) = \mathcal{L}^1(f^{(k-1)}) - \mathcal{L}^2(f^{(k-1)}) \text{ with } k = 1, \ldots, K.$$

### DeC Theorem

- $\mathcal{L}^1$ coercive
- $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz

DeC converges and $\min(K, M + 1)$ is the order of accuracy.

[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?

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

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

### DeC Theorem

- $\mathcal{L}^1$ coercive
- $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz

DeC converges and $\min(K, M + 1)$ is the order of accuracy.



[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?

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

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

$\mathcal{L}^1(f^{(k)}) = \mathcal{L}^1(f^{(k-1)}) - \mathcal{L}^2(f^{(k-1)})$ with $k = 1, \ldots, K.$

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

## DeC Theorem

- $\mathcal{L}^1$ coercive
- $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz

DeC converges and $\min(K, M + 1)$ is the order of accuracy.



[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?

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

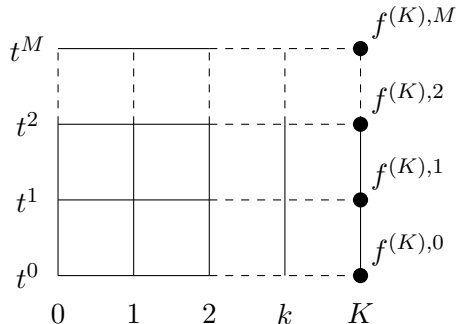$f^{m,(0)} := f(t^n), \quad m = 1, \ldots, M$

$\mathcal{L}^1(f^{(k)}) = \mathcal{L}^1(f^{(k-1)}) - \mathcal{L}^2(f^{(k-1)})$ with $k = 1, \ldots, K.$

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

### DeC Theorem

- $\mathcal{L}^1$ coercive
- $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz

DeC converges and $\min(K, M + 1)$ is the order of accuracy.



[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?

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

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

### DeC Theorem

- $\mathcal{L}^1$ coercive
- $\mathcal{L}^1 - \mathcal{L}^2$ Lipschitz

DeC converges and $\min(K, M + 1)$ is the order of accuracy.

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

## DeC – Proof

### Proof.

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

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

$$\alpha_1 ||f^{(k+1)} - f^*|| \leq ||\mathcal{L}^1(f^{(k+1)}) - \mathcal{L}^1(f^*)|| =$$

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

$$\leq \alpha_2 \Delta ||f^{(k)} - f^*||.$$

$$||f^{(k+1)} - f^*|| \leq \left(\frac{\alpha_2}{\alpha_1}\Delta\right) ||f^{(k)} - f^*|| \leq \left(\frac{\alpha_2}{\alpha_1}\Delta\right)^{k+1} ||f^{(0)} - f^*||.$$

After $K$ iteration we have an error at most of $\left(\frac{\alpha_2}{\alpha_1}\Delta\right)^K ||f^{(0)} - f^*||.$ □

Explicit DeC can be rewritten into Explicit Runge Kutta stages with $(r-1)^2 + 1$ stages
(with a correction due to the lumping of the mass matrix)

|              | Runge Kutta          | Deferred Correction                  |
|--------------|----------------------|--------------------------------------|
| Coefficients | Specific $\forall$ order | General algorithm                |
| Stages       | $r \leq s < r^2$     | $s = (r-1)^2 + 1$ or $(r-1\|\|r)$    |
| Mass matrix  | Full                 | Lumped                               |

## Idea of proof[6]

We know that

- $\mathcal{L}^1 = 0$ is AP.

We can prove that

- $\mathcal{L}_u^1 - \mathcal{L}_u^2 = \mathcal{O}(\varepsilon) + \mathcal{O}(\Delta)$
- $\mathcal{L}_f^1 - \mathcal{L}_f^2 = \mathcal{O}\left(\frac{\Delta}{\varepsilon}\right) + \mathcal{O}(\Delta)$.

---

[6]R. Abgrall, and D.T.. High Order Asymptotic Preserving Deferred Correction Implicit-Explicit Schemes for Kinetic Models. SIAM Journal on Scientific Computing, 42(3):B816–B845, 2020.

$$u_t + u_x = 0, \quad x \in [0,1], \quad t \in [0,T], T = 0.12, \quad u_0(x) = e^{-80(x-0.4)^2},$$

outflow BC, $\lambda = 1.5$, $\varepsilon = 10^{-10}$, $\theta_1 = 1$, $\theta_2 = 5$ (derivative stabilization).



(a) Scalar 1D convergence

(b) Order varying relaxation parameter

Next simulations will be over the Euler equation

$$\begin{pmatrix} \rho \\ \rho v \\ E \end{pmatrix}_t + \begin{pmatrix} \rho v \\ \rho v^2 + p \\ (E + p)v \end{pmatrix}_x = 0, \qquad x \in [0,1],\, t \in [0,T] \tag{15}$$

$\rho$ is the density, $v$ the speed, $p$ the pressure and $E$ the total energy. The system is closed by the equation of state

$$E = \frac{p}{\gamma - 1} + \frac{1}{2}\rho v^2. \tag{16}$$

# Numerical tests: Sod shock test

$\gamma = 1.4$, $T = 0.16$, outflow BC, $\varepsilon = 10^{-9}$, $\lambda = 2$, CFL $= 0.2$.
For $\mathbb{B}^1$ $\theta_1 = 1$, for $\mathbb{B}^2$ $\theta_1 = 1$, $\theta_2 = 0.5$, for $\mathbb{B}^3$ $\theta_1 = 2.5$, $\theta_2 = 4$.

$$\rho_0 = \mathbb{1}_{[0,0.5]}(x) + 0.1\mathbb{1}_{[0.5,1]}(x), \quad v_0 = 0, \quad p_0 = \mathbb{1}_{[0,0.5]}(x) + 0.125\mathbb{1}_{[0.5,1]}(x).$$
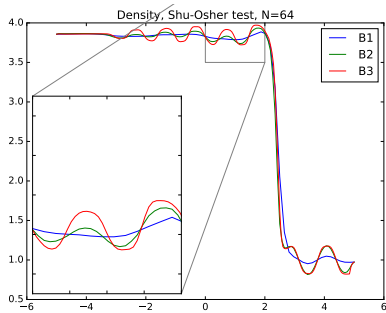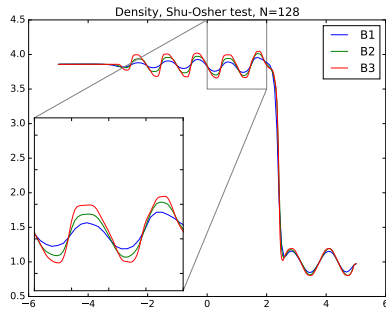


(a) $N = 64$      (b) $N = 256$

$\gamma = 1.4$, $T = 1.8$, outflow BC $\varepsilon = 10^{-9}$, $\lambda = 3$, CFL=0.1.
For $\mathbb{B}^1$ $\theta_1 = 0.5$, for $\mathbb{B}^2$ $\theta_1 = 0.8$, $\theta_2 = 1$, for $\mathbb{B}^3$ $\theta_1 = 3$, $\theta_2 = 1$.

$$\begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 3.857143 \\ 2.629369 \\ 10.333333 \end{pmatrix} x \in [-5, -4], \quad \begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 + 0.2\sin(5x) \\ 0 \\ 1 \end{pmatrix} \text{else.}$$
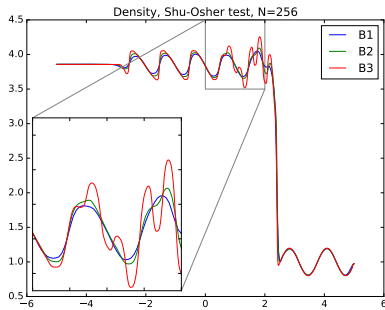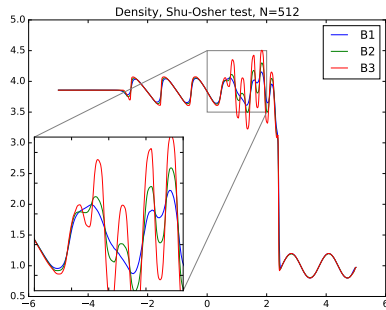


(c) $N = 64$

(d) $N = 128$

$\gamma = 1.4$, $T = 1.8$, outflow BC $\varepsilon = 10^{-9}$, $\lambda = 3$, CFL=0.1.
For $\mathbb{B}^1$ $\theta_1 = 0.5$, for $\mathbb{B}^2$ $\theta_1 = 0.8$, $\theta_2 = 1$, for $\mathbb{B}^3$ $\theta_1 = 3$, $\theta_2 = 1$.

$$\begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 3.857143 \\ 2.629369 \\ 10.333333 \end{pmatrix} x \in [-5, -4], \quad \begin{pmatrix} \rho_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 + 0.2\sin(5x) \\ 0 \\ 1 \end{pmatrix} \text{else.}$$



(e) $N = 256$

(f) $N = 512$

Euler equation in 2D domain

$$\partial_t U(\mathbf{x}, t) + \partial_x f(U(\mathbf{x}, t)) + \partial_y g(U(\mathbf{x}, t)) = 0, \, \mathbf{x} = (x, y) \in \Omega \subset \mathbb{R}^2,$$

$$U = \begin{pmatrix} \rho \\ \rho u \\ \rho v \\ E \end{pmatrix}, \quad f(U) = \begin{pmatrix} \rho u \\ \rho u^2 + p \\ \rho u v \\ u(E + p) \end{pmatrix}, \quad g(U) = \begin{pmatrix} \rho v \\ \rho u v \\ \rho v^2 + p \\ v(E + p) \end{pmatrix} \tag{17}$$

$\rho$ is the density, $u$ is the speed in $x$ direction, $v$ is the speed in $y$ direction, $E$ the total energy and $p$ the pressure.
The closing EOS is:

$$p = (\gamma - 1)\Big(E - \frac{1}{2}\rho(u^2 + v^2)\Big). \tag{18}$$

Initial conditions and solution for all $t \in [0, \infty)$ are

$$
\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} \left( 1 - \frac{\gamma-1}{\gamma} \frac{1}{2} \left( \frac{5}{2\pi} \right)^2 e^{\frac{1-r^2}{2}} \right)^{\frac{1}{\gamma-1}} \\ \frac{5}{2\pi}(-y)e^{\frac{1-r^2}{2}} \\ \frac{5}{2\pi}(x)e^{\frac{1-r^2}{2}} \\ \rho_0^\gamma \end{pmatrix}.
$$

Here $r^2 = x^2 + y^2$, the boundary conditions are outflow and $T = 1$.
$\gamma = 1.4$, $\varepsilon = 10^{-9}$, $\lambda = 1.4$ and CFL = 0.1.
For $\mathbb{B}^1$ $\theta_1 = 0.1$, for $\mathbb{B}^2$ $\theta_1 = 0.01$, $\theta_2 = 0$, for $\mathbb{B}^3$ $\theta_1 = 0.001$, $\theta_2 = 0$.

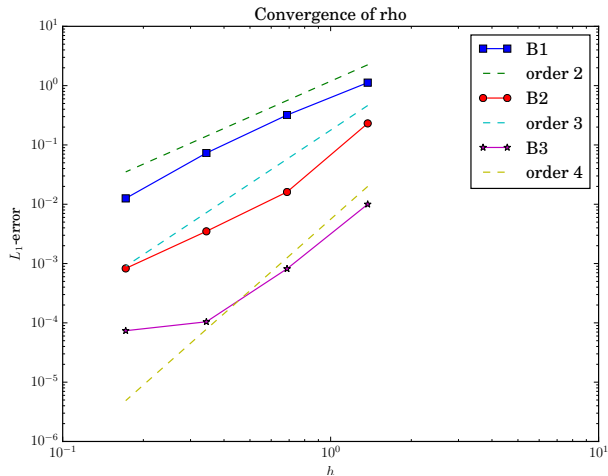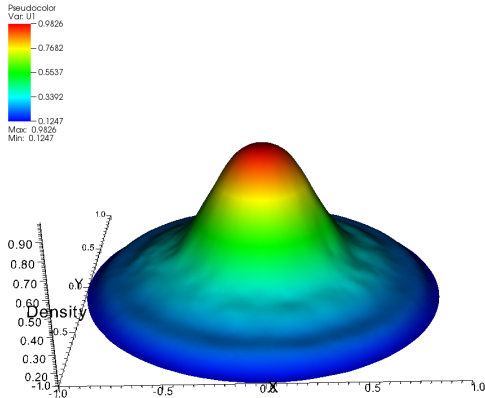Figure: 2D convergence

Initial conditions are

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 0 \\ 1 \end{pmatrix} \text{ if } r < \frac{1}{2}, \qquad \begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 0.125 \\ 0 \\ 0 \\ 0.1 \end{pmatrix} \text{ if } r \geq \frac{1}{2}.$$
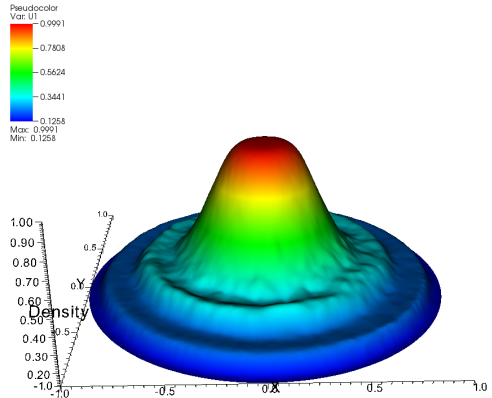
Here $r^2 = x^2 + y^2$, $\gamma = 1.4$, $\varepsilon = 10^{-9}$, $\lambda = 1.4$, CFL = 0.1, $T = 0.25$ and outflow boundary conditions.
For $\mathbb{B}^1$ $\theta_1 = 0.1$, for $\mathbb{B}^2$ $\theta_1 = 0.1$, $\theta_2 = 0.0001$, for $\mathbb{B}^3$ $\theta_1 = 0.01$, $\theta_2 = 0.0001$.
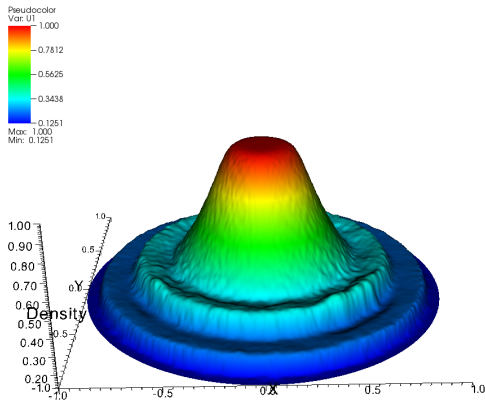
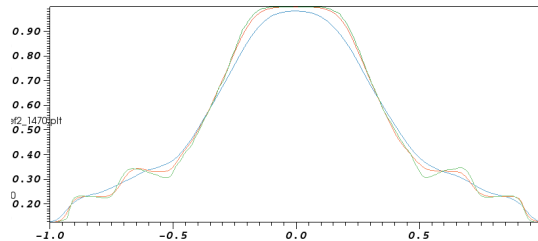# Numerical tests 2D: Sod shock test



(a) $\mathbb{B}^1, N = 13548$

(b) $\mathbb{B}^2, N = 13548$

(c) $\mathbb{B}^3, N = 13548$

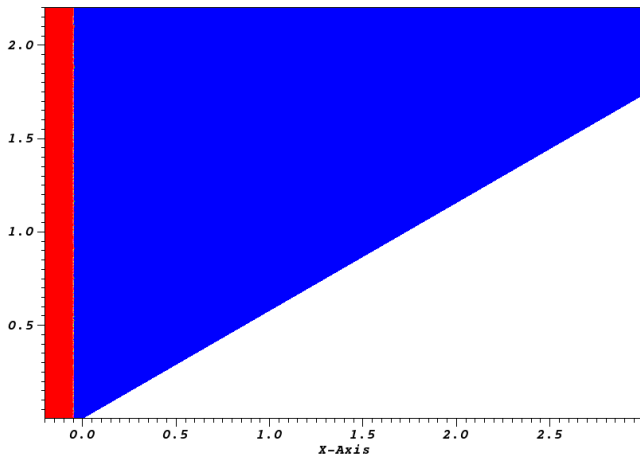(d) Slices of $\mathbb{B}^1$ (blue), $\mathbb{B}^2$ (red) and $\mathbb{B}^3$ (green), $N = 13548$

Double mach reflection test: initial conditions

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 8 \\ 8.25 \\ 0 \\ 116.5 \end{pmatrix} \text{ if } x \leq -0.05$$

$$\begin{pmatrix} \rho_0 \\ u_0 \\ v_0 \\ p_0 \end{pmatrix} = \begin{pmatrix} 1.4 \\ 0 \\ 0 \\ 1 \end{pmatrix} \text{ if } x > -0.05.$$

$T = 0.2$, $\varepsilon = 10^{-9}$, $\lambda = 15$, CFL $= 0.1$, $N = 19248$ triangular elements.
For $\mathbb{B}^1$ $\theta_1 = 0.1$, for $\mathbb{B}^2$ $\theta_1 = 0.01$, $\theta_2 = 0.0001$, for $\mathbb{B}^3$ $\theta_1 = 0.005$, $\theta_2 = 0.0001$.
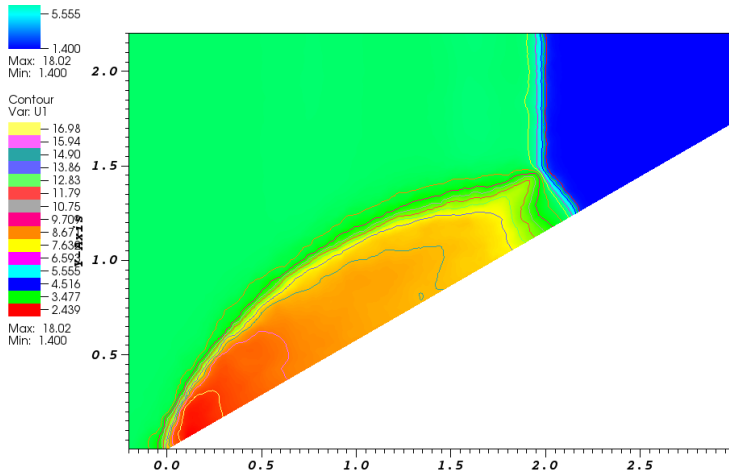
Figure: Density of DMR test $\mathbb{B}^1$

Figure: Density of DMR test $\mathbb{B}^2$

Figure: Density of DMR test $\mathbb{B}^3$

# Outline

# Shallow water equations

Modify the kinetic relaxation models by D. Aregba-Driollet and R. Natalini
Hyperbolic limit equation is

$$u_t + \sum_{d=1}^{D} \partial_{x_d} A_d(u) + S(u) = 0, \quad u : \Omega \to \mathbb{R}^K$$

$$\begin{cases} h_t + (hv)_x = 0 \\ (hv)_t + (hv^2 + \frac{g}{2}h^2)_x + ghb_x = 0 \end{cases}$$

Relaxation system

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right), \quad f^\varepsilon : \Omega \to \mathbb{R}^L$$

$$Pf^\varepsilon \to u, \quad P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

# Shallow water equations

Modify the kinetic relaxation models by D. Aregba-Driollet and R. Natalini
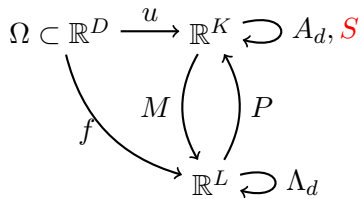Hyperbolic limit equation is

$$u_t + \sum_{d=1}^{D} \partial_{x_d} A_d(u) + S(u) = 0, \quad u : \Omega \to \mathbb{R}^K$$

$$\begin{cases} h_t + (hv)_x = 0 \\ (hv)_t + (hv^2 + \frac{g}{2}(h^2 - b^2))_x + g(h+b)b_x = 0 \end{cases}$$

Relaxation system

$$\Omega \subset \mathbb{R}^D \xrightarrow{u} \mathbb{R}^K \circlearrowright A_d, S$$

$$f \searrow \quad M \circlearrowright \quad P$$

$$\mathbb{R}^L \circlearrowright \Lambda_d$$

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right), \quad f^\varepsilon : \Omega \to \mathbb{R}^L$$

$$Pf^\varepsilon \to u, \quad P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u).$$

# Shallow water equations

Modify the kinetic relaxation models by D. Aregba-Driollet and R. Natalini
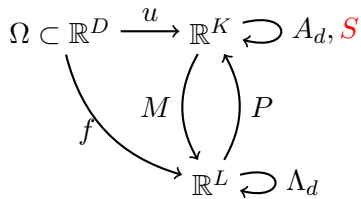Hyperbolic limit equation is

$$u_t + \sum_{d=1}^{D} \partial_{x_d} A_d(u) + S(u) = 0, \quad u : \Omega \to \mathbb{R}^K$$

$$\begin{cases} h_t + (hv)_x = 0 \\ (hv)_t + (hv^2 + \frac{g}{2}(h^2 - b^2))_x + g(h+b)b_x = 0 \end{cases}$$



Relaxation system

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon + \tilde{S}(f) = \frac{1}{\varepsilon}\left(M(Pf^\varepsilon) - f^\varepsilon\right), \quad f^\varepsilon : \Omega \to \mathbb{R}^L, \quad \tilde{S}(f) := \begin{pmatrix} S(f_1) \\ \cdots \\ S(f_N) \end{pmatrix},$$

$$Pf^\varepsilon \to u, \quad P(M(u)) = u, \quad P\Lambda_d M(u) = A_d(u), \quad P\tilde{S}(f) = S(Pf), \quad P\Lambda_d \tilde{S}(f) = S(P\Lambda f).$$

## Other properties

- Asymptotic preserving: Chapman–Enskog
- Well balancedness: lake at rest steady state preservation
  - Choice of a different form of the SW equation, so that the discretizations of the flux and the source match when $v = 0$
- Depth non-negativity
  - Wet and dry elements
  - Hybrid elements -> Modify the bathymetry to have positive DoFs

Figure: Subcritical flow: convergence for $\eta^\varepsilon = h^\varepsilon + b$ and $h^\varepsilon v^\varepsilon$

# Simulations: lake at rest



Figure: Lake at rest with immersed bump test: $\eta^\varepsilon$ and $v^\epsilon$ with $N = 25$

Figure: Lake at rest in parabola test: $\eta^\varepsilon$ and $v^\varepsilon$ with $N = 32$

# Simulations: Thucker Oscillations



Figure: Thacker oscillations in parabola test: $\eta^\varepsilon$ and $h^\varepsilon v^\varepsilon$ with $N = 100$

# Outline

## Conclusion and perspective

Conclusions

- Asymptotic preserving
- IMEX
- Residual Distribution
- Deferred Correction
- Idea for SW: well–balanced, wet/dry, nonnegative water height

Perspective

- Multiphase flows
- MOOD
- Entropy stability

# IMEX DeC RD – Bibliography

1. R. Abgrall, and D.T.. High Order Asymptotic Preserving Deferred Correction Implicit-Explicit Schemes for Kinetic Models. SIAM Journal on Scientific Computing, 42(3):B816–B845, 2020.

2. D. Aregba-Driollet and R. Natalini. Discrete kinetic schemes for multidimensional systems of conservation laws. SIAM J. Numer. Anal., 37(6):1973–2004, 2000.

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

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, 2017.

5. M. Ricchiuto, and A. Bollermann. Stabilized residual distribution for shallow water simulations. Journal of Computational Physics, 228(4):1071–1115, 2009.

Thank you for the attention!

Consider $M = 1$, $K = 2$.

$$\mathcal{L}^1(U^{(1)}, U^n) = 0. \tag{19}$$

$$\begin{cases} u_\sigma^{(1),n+1} = u_\sigma^n - \frac{\Delta t}{C_\sigma} \sum_{K|\sigma \in K} P\phi_\sigma^K(f^n) \\ f_\sigma^{(1),n+1} = \frac{\Delta t}{\varepsilon + \Delta t} M(u_\sigma^{(1),n+1}) + \frac{\varepsilon}{\Delta t + \varepsilon} f_\sigma^n - \frac{\varepsilon \Delta t}{C_\sigma(\Delta t + \varepsilon)} \sum_{K|\sigma \in K} \Phi_\sigma^K(f^n) \end{cases} \tag{20}$$

where $C_\sigma = \sum_{K|\sigma \in K} \int_K \varphi_\sigma(x) dx$.

## DeC – Example order 2 – Kinetic model

Consider $M = 1$, $K = 2$.

$$\mathcal{L}^1(U^{(2)}, U^n) = \mathcal{L}^1(U^{(1)}, U^n) - \mathcal{L}^2(U^{(1)}, U^n). \tag{21}$$

$$\begin{cases} u_\sigma^{(2),n+1} = u_\sigma^{(1),n+1} - \sum_{K|\sigma \in K} \int_K \varphi_\sigma (u^{(1),n} - u^n) + \\ \qquad\qquad - \frac{\Delta t}{C_\sigma} \sum_{K|\sigma \in K} P\left(\frac{1}{2}\phi_\sigma^K(f^n) + \frac{1}{2}\phi_\sigma^K(f^{(1),n+1})\right) \\ \\ f_\sigma^{(2),n+1} = f^{(1),n+1} + \frac{\Delta t}{\varepsilon + \Delta t}(M(u_\sigma^{(2),n+1}) - M(u_\sigma^{(1),n+1})) + \\ \qquad\qquad + \frac{\varepsilon}{\Delta t + \varepsilon} \sum_{K|\sigma \in K} \int_K \varphi_\sigma (f^{(1),n+1} - f^n) + \\ \qquad\qquad - \frac{\varepsilon \Delta t}{C_\sigma(\Delta t + \varepsilon)} \sum_{K|\sigma \in K} \frac{\Phi_\sigma^K(f^{(1),n+1}) + \Phi_\sigma^K(f^n)}{2} + \\ \qquad\qquad + \frac{\Delta t}{\Delta t + \varepsilon} \sum_{K|\sigma \in K} \int_K \varphi_\sigma \frac{M(u^{(1),n+1}) + M(u^n) - f^{(1),n+1} - f^n}{2} \end{cases} \tag{22}$$

where $C_\sigma = \sum_{K|\sigma \in K} \int_K \varphi_\sigma(x) dx$.

# Whitham's subcharacteristic condition

$$f_t^\varepsilon + \sum_{d=1}^{D} \Lambda_d \partial_{x_d} f^\varepsilon = \frac{1}{\varepsilon} \left( M(Pf^\varepsilon) - f^\varepsilon \right), \qquad f^\varepsilon : \Omega \to \mathbb{R}^L$$

If we call $u^\varepsilon = Pf^\varepsilon$, $v_d^\varepsilon = P\Lambda_d f^\varepsilon$ we have from (43) that

$$\begin{cases} \partial_t u^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} v_j^\varepsilon = 0 \\ \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} (P\Lambda_j \Lambda_d f^\varepsilon) = \frac{1}{\varepsilon} (A_d(u^\varepsilon) - v_d^\varepsilon) \end{cases}.$$

If we do a Taylor expansion in $\varepsilon$ we get

$$v_d^\varepsilon = A_d(u^\varepsilon) - \varepsilon \left( \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} (P\Lambda_d \Lambda_j f^\varepsilon) \right) \tag{23}$$

$$= A_d(u^\varepsilon) - \varepsilon \left( \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j} (P\Lambda_d \Lambda_j M(u^\varepsilon)) \right) + \mathcal{O}(\varepsilon^2). \tag{24}$$

# Whitham's condition

$$\partial_t u^\varepsilon + \sum_{d=1}^{D} \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^{D} \partial_{x_d} \left( \partial_t v_d^\varepsilon + \sum_{j=1}^{D} \partial_{x_j}(P\Lambda_d\Lambda_j M(u^\varepsilon)) \right) + \mathcal{O}(\varepsilon^2)$$

$$\partial_t u^\varepsilon + \sum_{d=1}^{D} \partial_{x_d} A_d(u^\varepsilon) = \varepsilon \sum_{d=1}^{D} \partial_{x_d} \left( \sum_{j=1}^{D} B_{dj}(u^\varepsilon)\partial_{x_j} u^\varepsilon \right) + \mathcal{O}(\varepsilon^2).$$

For this case, the Whitham's subcharacteristic condition[7] becomes

$$B_{jd} := P\Lambda_d\Lambda_j M'(u) - A_d'(u)A_j'(u), \qquad \sum_{j,d=1}^{D} (B_{dj}\xi_j, \xi_d) \geq 0.$$

---

[7]**natalini**.

## Problems: convection parameter

How to set the convection parameter automatically?
To verify Whitham's subcharacteristic condition we have to

$$B_{jd} := P\Lambda_d\Lambda_j M'(u) - A_d'(u)A_j'(u), \qquad \sum_{j,d=1}^{D} (B_{dj}\xi_j, \xi_d) \geq 0.$$

In DRM for 2D systems, we have:

$$\Lambda_1 = \begin{pmatrix} -\lambda I_K & 0_K & 0_K \\ 0_K & 0_K & 0_K \\ 0_K & 0_K & \lambda I_K \end{pmatrix}, \quad \Lambda_2 = \begin{pmatrix} 0_K & 0_K & 0_K \\ 0_K & -\lambda I_K & 0_K \\ 0_K & 0_K & \lambda I_K \end{pmatrix}$$

$$P\Lambda_1 = (-\lambda I_K, 0_K, \lambda I_K), \quad P\Lambda_2 = (0_K, -\lambda I_K, \lambda I_K)$$

$$P\Lambda_1\Lambda_1 = (\lambda^2 I_K, 0_K, \lambda^2 I_K), \quad P\Lambda_2\Lambda_2 = (0_K, \lambda^2 I_K, \lambda^2 I_K)$$

$$P\Lambda_1\Lambda_2 = P\Lambda_2\Lambda_1 = (0_K, 0_K, \lambda^2 I_K)$$

## Problems: convection parameter

Moreover we now that

$$
\mathbb{R}^{(K, K \cdot N)} \ni M'(u) =
$$

$$
= \begin{pmatrix} \frac{u}{3} + \frac{1}{3\lambda}(-2A_1 + A_2) \\ \frac{u}{3} + \frac{1}{3\lambda}(A_1 - 2A_2) \\ \frac{u}{3} + \frac{1}{3\lambda}(A_1 + A_2) \end{pmatrix}' = \frac{1}{3} \begin{pmatrix} I_K + \frac{1}{\lambda}(-2A_1' + A_2') \\ I_K + \frac{1}{\lambda}(A_1' - 2A_2') \\ I_K + \frac{1}{\lambda}(A_1' + A_2') \end{pmatrix}.
$$

So, if we compute the $B$ matrices we get

$$
B_{11} = \frac{2}{3}\lambda^2 I_K + \lambda(\frac{2}{3}A_2' - \frac{1}{3}A_1') - A_1' A_1'^T
$$

$$
B_{12/21} = \frac{1}{3}\lambda^2 I_K + \lambda(\frac{1}{3}A_2' + \frac{1}{3}A_1') - A_{1/2}' A_{2/1}'^T
$$

$$
B_{22} = \frac{2}{3}\lambda^2 I_K + \lambda(\frac{2}{3}A_1' - \frac{1}{3}A_2') - A_2' A_2'^T
$$

## Problems: convection parameter

Then, if we restart from the following condition

$$\sum_{i,j=1}^{2} \langle B_{ij}\xi_i, \xi_j \rangle \geq 0 \qquad \forall \xi_j \in \mathbb{R}^K,$$

Different from scalar case $K = 1$. Scalar case:

$$\sum_{i,j=1}^{2} \langle B_{ij}\xi_i, \xi_j \rangle \geq 0 \qquad \forall \xi_j \in \mathbb{R},$$

you can get something solvable, but in our case, what we get is:

$$\frac{2}{3}\sum_{i,j=1}^{2} \langle \xi_i, \xi_j \rangle \lambda^2 + \frac{\lambda}{3}\big(\langle (2A_2' - A_1')\xi_1, \xi_1 \rangle +$$
$$+ \langle (-A_2' + 2A_1')\xi_2, \xi_2 \rangle + \langle (A_2' + A_1' + (A_2' + A_1')^T)\xi_1, \xi_2 \rangle \big) +$$
$$+ \sum_{i,j=1}^{2} \langle A_i'A_j'^T\xi_i, \xi_j \rangle \geq 0, \qquad \forall \xi_1, \xi_2 \in \mathbb{R}^K.$$

## Problems: convection parameter

How they saw this was in the sense of

$$\underline{\xi}^T B \underline{\xi} \geq 0.$$

So doing spectral analysis, finding the eigenvalues of $B$ and imposing the positivity of both of them for *scalar* case. Finally, they got this condition from a 4th degree equation

$$\lambda \geq \max\left(-A_1' - A_2', 2A_1' - A_2', -A_1' + 2A_2'\right).$$

But for general case $B$ is a $2K \times 2K$ matrix and I have no clue how to find the $2K$ eigenvalues.

If we change the convection parameter from timestep to timestep, we get big oscillations.
Where should this come from?
Back to IMEX 3

## Residual distribution - Choice of the scheme

How to split into $\phi_\sigma^K \Rightarrow$ choice of the scheme. For example, we can rewrite SUPG in this way:

$$\phi_\sigma^K(U_h) = \int_K \varphi_\sigma(\nabla \cdot A(U_h) - S(U_h))dx+ \tag{25}$$

$$+h_K \int_K (\nabla \cdot A(U_h) \cdot \nabla \cdot \varphi_\sigma)\, \tau\, (\nabla \cdot A(U_h) \cdot \nabla \cdot U_h)\,. \tag{26}$$

Furthermore, we can write the Galerkin FEM scheme with jump stabilization by **burman**:

$$\phi_\sigma^K = \int_K \varphi_\sigma(\nabla \cdot A(U_h) - S(U_h))dx + \sum_{e|\text{edge of } K} \theta h_e^2 \int_e [\nabla U_h] \cdot [\nabla \varphi_\sigma]d\Gamma, \tag{27}$$

## Residual Distribution - Choice of the scheme

$$\phi_\sigma^{K,LxF}(U_h) = \int_K \varphi_\sigma \left(\nabla \cdot A(U_h) - S(U_h)\right) dx + \alpha_K(U_\sigma - \overline{U}_h^K), \tag{28}$$

where $\overline{U}_h^K$ is the average of $U_h$ over the cell $K$ and $\alpha_K$ is defined as

$$\alpha_K = \max_{e \text{ edge } \in K} \left(\rho_S \left(\nabla A(U_h) \cdot \mathbf{n}_e\right)\right), \tag{29}$$

$\rho_S$ is the spectral radius.
For monotonicity near strong discontinuities, PSI limiter:

$$\beta_\sigma^K(U_h) = \max\left(\frac{\Phi_\sigma^{K,LxF}}{\Phi^K}, 0\right) \left(\sum_{j \in K} \max\left(\frac{\Phi_j^{K,LxF}}{\Phi^K}, 0\right)\right)^{-1} \tag{30}$$

Blending between LxF and PSI:

$$\phi_\sigma^{*,K} = (1 - \Theta)\beta_\sigma^K \phi_\sigma^K + \Theta\Phi_\sigma^{K,LxF},$$
$$\Theta = \frac{|\Phi^K|}{\sum_{j\in K}|\Phi_j^{K,LxF}|}. \tag{31}$$

Nodal residual is finally given by

$$\phi_\sigma^K = \phi_\sigma^{*,K} + \sum_{e|\text{edge of } K} \theta h_e^2 \int_e [\nabla U_h] \cdot [\nabla \varphi_\sigma] d\Gamma. \tag{32}$$

## DeC – Proof

### Proof.

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

$$
\begin{aligned}
\mathcal{L}^1(U^{(k+1)}) - \mathcal{L}^1(U^*) &= \left( \mathcal{L}^1(U^{(k)}) - \mathcal{L}^2(U^{(k)}) \right) - \left( \mathcal{L}^1(U^*) - \mathcal{L}^2(U^*) \right) \\
&= \left( \mathcal{L}^1(U^{(k)}) - \mathcal{L}^1(U^*) \right) - \left( \mathcal{L}^2(U^{(k)}) - \mathcal{L}^2(U^*) \right) \\
\alpha_1 ||U^{(k+1)} - U^*|| &\leq ||\mathcal{L}^1(U^{(k+1)}) - \mathcal{L}^1(U^*)|| = \\
&= ||\mathcal{L}^1(U^{(k)}) - \mathcal{L}^2(U^{(k)}) - (\mathcal{L}^1(U^*) - \mathcal{L}^2(U^*))|| \leq \\
&\leq \alpha_2 \Delta ||U^{(k)} - U^*||. \\
||U^{(k+1)} - U^*|| &\leq \left( \frac{\alpha_2}{\alpha_1} \Delta \right) ||U^{(k)} - U^*|| \leq \left( \frac{\alpha_2}{\alpha_1} \Delta \right)^{k+1} ||U^{(0)} - U^*||.
\end{aligned}
$$

After $K$ iteration we have an error at most of $\eta^K \cdot ||U^{(0)} - U^*||$. $\qquad \square$